

LH_{RS}*—A Highly-Available Scalable Distributed Data Structure

WITOLD LITWIN and RIM MOUSSA

Université Paris Dauphine

and

THOMAS SCHWARZ, S. J.

Santa Clara University

LH_{RS}* is a high-availability scalable distributed data structure (SDDS). An LH_{RS}* file is hash partitioned over the distributed RAM of a multicomputer, for example, a network of PCs, and supports the unavailability of any $k \geq 1$ of its server nodes. The value of k transparently grows with the file to offset the reliability decline. Only the number of the storage nodes potentially limits the file growth. The high-availability management uses a novel parity calculus that we have developed, based on Reed-Salomon erasure correcting coding. The resulting parity storage overhead is about the lowest possible. The parity encoding and decoding are faster than for any other candidate coding we are aware of. We present our scheme and its performance analysis, including experiments with a prototype implementation on Wintel PCs. The capabilities of LH_{RS}* offer new perspectives to data intensive applications, including the emerging ones of grids and of P2P computing.

Categories and Subject Descriptors: E.1 [Data]: Data Structures—*Distributed data structures*; D.4.3 [Operating Systems]: File Systems Management—*Distributed file systems*; D.4.5 [Operating Systems]: Reliability—*Fault-tolerance*; H.2.2 [Database Management]: Physical Design—*Access methods; Recovery and restart*

General Terms: Algorithms, Reliability

Additional Key Words and Phrases: Scalable distributed data structure, linear hashing, high-availability, physical database design, P2P, grid computing

Motto: Here is Edward Bear, coming downstairs now, bump, bump, bump, on the back of his head, behind Christopher Robin. It is, as far as he knows, the only way of coming downstairs, but sometimes he feels that there really is another way, if only he could stop bumping for a moment and think of it. And then he feels that perhaps there isn't.

This research was partially supported by a generous gift from Microsoft Research, Europe and EEC-ICONS project no. IST-2001-32429.

Authors' addresses: W. Litwin, R. Moussa, Université Paris 9 (Dauphine), Pl. du Mal. Du Lattre, Paris 75016, France; email: {Witold.Litwin,Rim.Moussa}@dauphine.fr; T. Schwarz, S. J., Department of Computer Engineering, Santa Clara University, 500 El Camino Real, Santa Clara, CA 95053-0566; email: tjschwarz@scu.edu.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or direct commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 1515 Broadway, New York, NY 10036 USA, fax: +1 (212) 869-0481, or permissions@acm.org.

© 2005 ACM 0362-5915/05/0900-0769 \$5.00

Winnie-the-Pooh. By A. A. Milne, with decorations by E. H. Shepard. Methuen and Co, London (publ.)

1. INTRODUCTION

Shared-nothing configurations of computers connected by a high-speed link, often called *multicomputers*, allow for high aggregate performance. These systems gained in popularity with the emergence of grid computing and P2P applications. They need new data structures that scale well with the number of components [Com. ACM 1997]. Scalable Distributed Data Structures (SDDS) aim to fulfill this need [Litwin et al. 1993; SDDS-bibliography]. An SDDS file is stored at multiple nodes provided by *SDDS servers*. As the file grows, so does the number of servers on which it resides. The SDDS addressing scheme has no centralized components. This allows for operation speeds independent of the file size. They provide for hash, range or m-d partitioned files of records identified by a primary or by multiple keys. See SDDS [bibliography] for a partial list of references. A prototype system, SDDS 2000, for Wintel PCs, is freely available for a non-commercial use [CERIA <http://ceria.dauphine.fr>].

Among the best-known SDDS schemes is the LH* scheme [Litwin et al. 1993, 1996; Karlson et al. 1996; Breitbart et al. 1996; Bertino et al. 1999; Knuth 1993; Ramakrishnan 1999]. LH* creates scalable, distributed, hash-partitioned files. Each server stores the records in a bucket. The buckets split when the file grows. The splits follow the *linear hashing* (LH) principles [Litwin 1980, 1994]. Buckets are stored for fast access in distributed RAM, otherwise they can be on disks. Only the maximum possible number of server nodes limits the file size. A search or an insert of a record in an LH* file can be hundreds of times faster than a disk access [Bennour et al. 2000; Bennour 2002].

At times, an LH* server can become unavailable. It may fail as the result of a software or hardware failure. It may also stop the service for good or for an unacceptably long time, a frequent case in P2P applications. Either way, access to data becomes impossible. The situation may not be acceptable for an application, limiting the utility of the LH* scheme. Data unavailability can be very costly, [www.contingencyplanningresearch.com/cod.htm/1996]. An unavailable financial database may easily cost the owner \$10K–\$27K per minute [Bartalos 1999].

A file might suffer from the unavailability of several of its servers. We say that it is *k-available*, if all data remain available despite the unavailability of any *k* servers. The information-theoretical minimum storage overhead for *k*-availability of *m* data servers is k/m [Hellerstein et al. 1994]. It requires *k* additional, so-called parity symbols (records, buckets . . .) per *m* data symbols (records, buckets . . .). Decoding *k* unavailable symbols requires access to *m* available symbols of the total of $m + k$. Large values for *m* seem impractical. A reasonable approach to limit *m* is to partition a data file into groups consisting of at most *m* nodes (buckets) per group, with independent parity calculus.

For files on a few servers, 1-availability usually suffices. The parity calculus can then be the fastest known, using only XORing, as in RAID-5. The probability of a server unavailability increases however with the file size. The file *reliability*,

which is the probability that all the data are available for the application, necessarily declines. At one point we need 2-availability despite the increased storage overhead. Likewise, as the file continues to grow, at some point we need 3-availability, despite further increase to the storage overhead. In fact, for any fixed k , the probability of k -unavailability increases with the file size. For largely scaling files, we need the *scalable availability*, adjusting k to the file size [Litwin et al. 1998].

Below, we present an efficient scalable availability scheme we called LH_{RS}*. It structures the LH* data buckets into groups of size m , providing each with $K \geq 1$ parity buckets. The values of m and of K are file parameters that can be adjusted dynamically. We call K the *intended availability level*. A group is typically K -available. Some can be $(K-1)$ -available with respect to the current K if the file just increased it by one. The new level diffuses to the groups progressively with the splits. Changes to K value are transparent for the application.

The scheme's theory originates in Litwin and Schwarz [2000]. Since then, we have sped up the parity calculus, without compromising the storage overhead. We also have progressively optimized implementation issues, especially the efficiency of the communication architecture. We have realized a prototype for the multi-computer of Wintel PCs on a typical LAN [Litwin et al. 2004a,b]. We have analytically and experimentally determined various performance factors. The work has validated our design choices. Below we report on all these issues.

Our current parity calculus is a novel erasure correction scheme, using Reed-Solomon (RS) coding. Our storage overhead for k -availability is close to the optimal one of k/m . To our best knowledge, our schemes offer the fastest parity generation (encoding) for our needs. We have also optimized the erasure correction (decoding), although with a lesser priority, expecting it to be much less frequent, hopefully.

More specifically, we recall that an RS code uses a Galois Field (GF). The addition in a GF amounts to XORing only. Multiplication is necessarily slower [MacWilliams and Sloane 1997]. In departure from [Litwin and Schwarz 2000], we determined the $GF(2^{16})$ and $GF(2^8)$ to be more appropriate. Our scheme uses now only XORing to generate the first parity symbol (record, bucket, ...). We also only need XORing for the recovery of a single data bucket. In addition, changes to a data record in the first bucket of any group result in XORing only at the parity buckets, further speeding up $1/m$ of our encoding operations. Moreover, we have accelerated the k -parity encoding and decoding, by introducing the *logarithmic* parity and decoding matrices. All this makes our parity encoding scheme the fastest known, under the minimal group storage overhead constraint, to our best knowledge of the k -erasure correction codes. Besides, our high-availability features are transparent to the application and do not affect the speed of searches and scans in an LH_{RS}* file. These perform as well as in an LH* file with the same data records and bucket size.

LH_{RS}* is the only high-availability SDDS scheme prototyped to the extent that we present. However, it is not the only one known. Some proposals use mirroring to achieve 1-availability [Litwin and Neimat 1996; Breitbart and Vingralek 1998; Vingralek et al. 1998]. Two schemes use only XORing to provide 1-availability [Litwin et al. 1997; Litwin and Risch 2001; Lindberg

1997]. Another XORing-only scheme LH_{SA}^* was the first to offer scalable availability [Litwin et al. 1999]. It can generate parity faster than LH_{RS}^* , but sometimes has greater storage overhead. We compare various schemes in the Related Work section.

To address the overall utility of the LH_{RS}^* scheme, we recall that the existing hash file schemes store files on local disk(s), or statically partition (parallelize) them over a cluster of nodes. The latter approach provides for larger files or better efficiency of the (non-key) scans. As we have mentioned, the disks may use underneath a hardware or software RAID scheme, 1-available RAID5 typically. This may adversely affect a record access performance, for example, by segmenting some records, or manipulating blocs much larger than records, as it works independently of the hash scheme. Partitioned files are sometimes replicated, for 1-availability typically given the incurred storage overhead cost.

The LH_{RS}^* scheme is “plug compatible” with such schemes. It offers the same functional capabilities: the key search, a scan, a record insert, update, and delete. The access performance should in contrast, typically improve by orders of magnitude, especially for files that can now entirely fit for processing into the potentially unlimited (distributed) RAM of a multicomputer, while they could not at a single node or a cluster. Our experiments with LH_{RS}^* files in distributed RAM showed the individual key search to be about 30 times faster than a single access to a local disk. A bulk search speeds up 200 times. Experiments with inserts lead to similar figures, for up to a 3-available file. The application may thus accelerate the query processing from, let’s say 30 minutes or 3 hours, to a single minute.

Likewise, the administrator may choose the availability level “on-demand” or may let it autonomously adjust to the file size. The storage overhead is always about the minimal possible. Both properties are uniquely attractive at present for larger and very large (Pbyte) files. Records are never segmented and data transfers from the storage are not bigger, as the high-availability is designed to work in accordance with the hash scheme in this way. Our experiments showed furthermore that the recovery from a k -unavailability should be fast. For instance, about 1.5 seconds sufficed to recover more than 10MB of data (100 000 records) within three unavailable data buckets. Next, the file may become very large, spreading dynamically over theoretically any number of nodes. The scaling transparency finally frees the application administrator from periodic file restructuring. Some are already eagerly waiting for this perspective to materialize [Ben-Gan and Moreau 2003].

Hash files are ubiquitous. Legions of their applications could benefit from the new capabilities of LH_{RS}^* . This is particularly true for the numerous applications of major DBMSs. DB2 and Oracle, which use single site or parallel hash files, and Postgres DBMS, which uses single site linear hash files. Other applications affected by the current technology include video servers, Web servers, high performance file servers, and dedicated mission-critical servers. In the latter category, the rapidly and constantly growing data of a well-known large-scale search engine may allegedly already need about 54.000 nodes [www.economist.com].

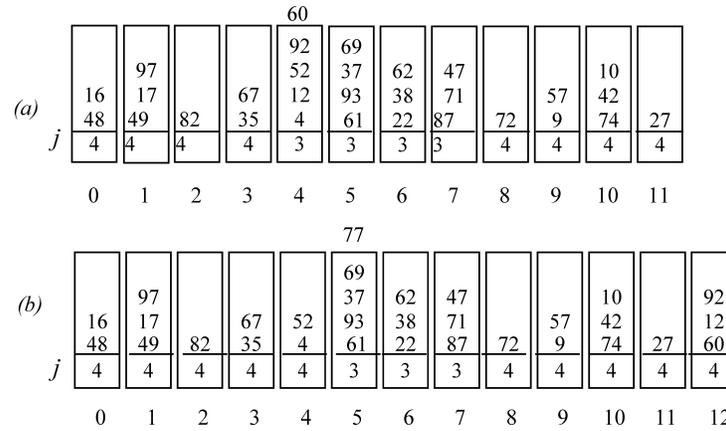


Fig. 1. LH_{RS}* file data buckets, before (a) and after (b) the insert of record 77.

The presently unrivaled scalability of LH_{RS}* files should also serve the emerging needs of data grids and of P2P applications, as we discuss in Section 6. In particular, a number of new applications target files larger than anything operationally workable till now, for example, the SkyServer project to cite just one [Gray et al. 2002]. On a different note, the sheer access speed to data in an LH_{RS}* file together with its ability to quickly scale, makes it an attractive tool for stream analysis. In this increasingly popular domain, streams often come simultaneously from multiple sensors, thus needing storage at the generation speed. Incidentally, the problem of collection of such streams from an airplane engine, at Santa Clara University mechanical lab, triggered the original LH* scheme idea.

We describe the general structure of an LH_{RS}* file in Section 2. Section 3 presents the parity calculus. We explain the LH_{RS}* file manipulations in Section 4. Section 5 deals with the performance analysis. Section 6 discusses related work. Section 7 concludes the study and proposes directions for future work. Appendix A shows our parity matrices for $GF(2^{16})$ and $GF(2^8)$. Appendix B sums up our terminology. We give additional details in Online Appendix C available in the ACM Digital Library regarding parity calculus, file operations, performance analysis and variants of the basic scheme, including a discussion of alternative erasure correcting codes.

2. FILE STRUCTURE

LH_{RS}* provides high availability to the LH* scheme [Litwin et al. 1993, 1996; Karlson et al. 1996; Litwin et al. 1999]. LH* itself is the scalable distributed generalization of Linear Hashing (LH) [Litwin 1980a,b]. An LH_{RS}* file stores *data* records in data buckets, numbered 0, 1, 2, . . . and *parity* records in separate parity buckets. Data records contain the application data. A data bucket has the *capacity* of $b \gg 1$ *primary* records. Additional records become *overflow* records. The application interacts with the data records as in an LH* file. Parity records only provide the high availability and are invisible to the application.

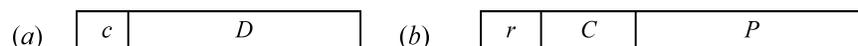


Fig. 2. LH_{RS}^* record structure: (a) data record, (b) parity record.

We store the data and parity buckets at the LH_{RS}^* server nodes. There is basically one bucket per server. The application does not address any server directly. It calls an LH_{RS}^* client component, usually at the application node. The clients address the servers for data search or storage over the network, transparently for the applications.

An LH_{RS}^* operation is in *normal* mode as long as it does not encounter an unavailable bucket. If it does, then it enters *degraded* mode. Below, we assume the normal mode unless otherwise stated. We first present the storage and addressing of the data records. We introduce the parity management afterwards.

2.1 Data Records

The storage and addressing of data records in an LH_{RS}^* file is the same as in an LH_{LH}^* file [Karlson et al. 1996]. We thus only briefly recall the related principles. Details are in Karlson et al. [1996], as well as in Litwin et al. [1996], Bennour et al. [2000], and Bennour [2002]. We follow up with the description of the LH_{RS}^* specific rules for *misdirected* requests.

2.1.1 Storage. A data record structure consists of a *key* identifying it and of other, non-key, data. Figure 2a, shows the record structure, with c denoting the key and D the nonkey field, usually much longer than the key field. The key determines the record location (the bucket number a) through the LH-function [Litwin 1980]:

(LH) $a := h_i(c)$; **if** $a < n$ **then** $a := h_{i+1}(c)$.

Here, h_i stands for a family of hash functions with specific properties. Usually, one uses $h_i = c \bmod 2^i$. The variables (i, n) are the *file state*, where $i = 0, 1 \dots$ stands for the *file level* and $n = 0, 1 \dots$ is the *split pointer*. The state, hence the LH-function itself, dynamically adjusts to the file size. The state determines the number N of data buckets in the file as $N = 2^i + n$. We call N *file extent*. Vice versa, the extent determines the file state. For every N , 2^i is the largest power of 2 smaller or equal to N .

The initial file state is $(0, 0)$ for the file extent $N = 1$. Every record then hashes to bucket 0. LH^* file adjusts to the scale-up by bucket splits. A low ratio of overflows result from typically and good access performance. If the file shrinks, buckets can merge, to prevent an underload. It appears that merges are not used in practice for file systems. Therefore, we forego further discussion.

An LH_{RS}^* component called the *coordinator* manages the splits (and merges). It may reside anywhere, but it should be practical to locate it at the node of bucket 0.¹ Both components however remain distinct: the unavailability of bucket 0 does not imply that of the coordinator and vice versa. The coordinator can be k -replicated for its own k -availability, having very little data on

¹For LH^* , there are also variants without the coordinator [Litwin et al. 1996], but we are not aware of attempts to make them highly available.

its own, as will appear. When an insert causes an overflow, the bucket informs the coordinator. The coordinator sends the split request to bucket n . It also appends to the file bucket $N = 2^i + n$, allotted at some server, according to some allocation scheme [Litwin et al. 1996]. Bucket n readdresses (rehashes) every record in it using h_{i+1} . A crucial property of LH-functions is that h_{i+1} only maps each key to n or to $2^i + n$. The records mapped to n , half of those in the bucket on the average, remain. The others move to a new bucket. Both buckets retain the value $i + 1$ as the *bucket level*. The initial level of bucket 0 is zero.

Each split increases n to $n + 1$. The resulting (LH) calculus with the new state conveniently takes the new bucket into account. The buckets split in this way in numerical (linear) order, until all the buckets in the file reach the level $i + 1$. This occurs when all the buckets 0 to $2^i - 1$ have split using h_{i+1} . The coordinator then increases the file level to $i := i + 1$ and moves the split pointer back to $n := 0$. The splitting process uses the new h_{i+1} from then on.

The internal structure of an LH_{RS}* data bucket is that of an LH_{LH}* bucket. Each bucket is basically a RAM LH file, as in Litwin [1980]. The internal buckets called pages are the chains of records; the overflow records are those beyond the b -th element. This result provides for efficient page splitting and fast local access. Perhaps more surprisingly, it also accelerates the LH_{RS}* splits, with respect to their naïve processing.

Example 1. Consider the LH_{RS}* data buckets in Figure 1a. The split pointer is $n = 4$, the file level is $i = 3$, the file extent is $N = 12 = 2^3 + 4$ buckets. The header of each bucket shows its level j . We only show the key fields c , assumed to be integers. The bucket capacity is $b = 4$. Bucket 4 has an overflow record. Its insertions triggered one of the previous splits. Most likely, one of those using h_4 that processed the file up to bucket 3 until then. We now insert record 77. According to (LH) it enters bucket 5. The bucket reports an overflow to the coordinator. This overflow creates bucket $N = 12$, Figure 1b, and requests bucket 4 to split. Bucket 4 has level $j = 3$, hence rehashes all its records using h_4 . The new address for a record can be only either 4 or 12. The bucket keeps every record hashed to 4. It sends every other record to bucket 12. In our case, two records remain, while three move. Both buckets get level $j = 4$. Bucket 4 reports the end of the split to the coordinator, which finally moves the split pointer n to $n + 1 = 5$. Now, bucket 4 no longer overflows, but even has room for three more primary records. Likewise, there is room for two more inserts in new bucket 12. The overflow at bucket 5 remains though, despite being at the origin of the split. Its resolution waits for coming splits. In our case, the next one is in fact at bucket 5 and will indeed resolve the overflow.

2.1.2 Addressing. As in LH* file, each LH_{RS}* client caches its private image of the file state, and applies (LH) to it. The coordinator would indeed become a hot spot if the clients should access it for the addressing. The coordinator does not push every file state update to the clients nor even to the servers, for the same reason. A split therefore makes every existing image outdated. The initial image of a new client is $(0, 0)$. A client with an outdated image may direct a key-based request towards an *incorrect* address, not the *correct* one given by

(LH) for the request. The addressee recognizes its status from the received key c and from its bucket level j . The correct address must indeed be equal to $h_j(c)$ that is its own.

The incorrectly addressed bucket should forward the request to the bucket that is possibly the correct one for c , and within the shortest extent N that could exist in the situation. This guarantees that the hop does not exceed any current extent. The guess can be based solely on j and the bucket's own address. Indeed, only the coordinator knows n . The first guess must be bucket $h_{j-1}(c)$. This address must be under $2^j - 1$, hence the guess is always safe, as the extent of the file with at least one bucket of level j must be beyond that bound. If this address is not that of the guessing bucket itself, then it cannot verify whether the level of the addressee is actually $j - 1$ or j . Hence, it has to resend the request there. Otherwise, the guess does not lead to any actual forwarding. The only remaining choice is $h_j(c)$. This must be the correct address. If the request is resent to bucket $h_{j-1}(c)$, and the address also happens to be incorrect, then the bucket level was indeed equal to j . The same calculus there resends then the request to the correct bucket $h_j(c)$.

An incorrect bucket receiving the request from the client resends it according to the above reasoning. It executes the resulting LH_{RS}^* *Forwarding Algorithm* shown in App. C in Online Appendix C available in the ACM Digital Library. The addressee may resend it once more, using the same algorithm. This must be the last hop. In other words, any client's request dealt with in this way must reach its correct bucket in at most two hops. This property is independent of the file extent. It is unique to LH^* based schemes up to now, and comparatively, among SDDSs, the most efficient known at present. Besides, the performance analysis for LH^* has shown that most of the key-based addressing requests in fact do not encounter any forwarding in practice, for a reasonable large bucket capacity b .²

The correct bucket sends an *Image Adjustment Message* (IAM) to the client. It contains essentially the bucket level and the address of the latest resending bucket.³ The client updates its image on this basis. It guesses the shortest file extent and the related state that could lead to the data it got. The algorithm avoids making the same addressing error twice. It also makes the client's image of the extent closer to the actual one. Numerically, App. C of Online Appendix C available in the ACM Digital Library, if j is the received level, and A is the address, the image evaluates to $(i' = j - 1, n' = A + 1)$, unless n' then turns out to be $2^{i'}$. The file state reevaluates then instead to $(i' = j, n' = 0)$. The rationale for the calculus is largely similar to that already outlined for the forwarding calculus at the bucket.

The IAMs may also refresh the client's image of the physical addresses of the servers. Each IAM brings to the client the addresses of all the buckets stretching

²Incidentally, these properties nicely fit D. Knuth's comments about amazing efficiency of randomized algorithms in Lecture 2 of *Things a Computer Scientist Rarely Talks About*, CSLI Publications, Stanford, 2001.

³This choice is actually a novelty for LH^* based schemes. Up to now they used the one in Litwin et al. [1996]. The new approach improves the image adequacy in many cases.

between the one initially addressed by the client and the correct one. A specific IAM we describe later also refreshes the client when the recovery displaces a bucket.

Example 2. In the file at Figure 1a, the key $c = 60$ is in bucket 4 according to (LH). Now, a client with the file state image equal to $(0, 0)$, hence with the file extent image equal to $N' = 1$, searches for this key. The client sends the request to bucket 0, according to its execution of (LH). Bucket 0 verifies whether it could be the correct bucket for record 60. Its bucket level is 4, hence we have $h_4(60) \neq 0$. The bucket is not the correct one, hence it needs to forward the request. The bucket guesses the new address as $h_3(60) = 4$. It is safely in the existing file extent, as the file with bucket 0 of level 4 must have at least 9 buckets, while h_3 may only address buckets $0 \dots 7$. The forwarding address does not point to bucket 0 itself, hence the bucket forwards the request to bucket 4. It includes in the message its own bucket level $j = 4$. The calculus at bucket 4 yields $h_3(60) = 4$. Hence it is the correct bucket. It therefore responds to the client with record 60 (if the request was a key search). It piggy-backs the IAM containing the address 0 and the received level $j = 4$. The client updates the file state image to $(3, 1)$. This amounts to a new extent image of $N' = 9$. This value is much closer to the actual one of $N = 12$ than the initial $N' = 1$. If the client repeats the search, it sends the request directly to bucket 4, avoiding the previous error. Moreover, initially any request from the client for any but the two records in bucket 0 would require the forwarding. Now, only the requests for the six records in buckets 9–11 would require it.

Consider now another client with image of $N = 1$ sending the request for record 60 to the file at Figure 1b. Bucket 0 again forwards the request, and its bucket level, to bucket 4. The calculus at bucket 4 now yields $h_4(60) = 12$. The request cannot be resent to bucket $h_3(60)$ since it is bucket 4 itself. Hence the bucket uses the guess of $h_4(60)$ and resends the request, with its own bucket level to bucket 12. This has to be the correct address. Bucket 12 handles the request and sends the IAM that the level of bucket 4 is 4. The client's new image is $(3, 5)$ hence it amounts to the image $N' = 13$. It is the perfect guess. An access to any record by our client now arrives directly at the correct bucket, until next split.

2.1.3 Misdirected Requests. LH_{RS}* restores an unavailable bucket on a spare server. It is typically a different server than the one with the unavailable bucket. Only the client and the buckets involved in the recovery are aware of the new location. Any other client or bucket can *misdirect* a request to the former server. A misdirected request could also follow a bucket merge. The LH_{RS}* file processes the misdirected requests basically as follows.

Any request to a data bucket carries the number of its intended bucket. A node that a request reaches is supposed basically to be an SDDS server. Other nodes are not supposed to react in any way. Like, *ipso facto*, an unavailable server. The request without reply enters the degraded mode that we discuss in Section 4. The server getting a request verifies that it carries the intended bucket. If so, it acts as described above. Otherwise, it forwards the request to

the coordinator. This one resends the request to the correct bucket. An IAM informs the sender of the location change.

The scheme is valid provided that a misdirected request never finds the bucket at a former location to be available despite its recovery elsewhere. If it could happen, the client could unknowingly get data that is not up-to-date. The assumptions for the LH_{RS}^* design that we now outline prevent this eventuality. They are rooted in the popular fail-stop model and fit well the local network environment we mainly target.

We thus presume that when a client or server finds a bucket unavailable, the unavailability is permanent and of the entire bucket. Next, we suppose that the communications are (fully) reliable. Finally, we presume that the server is (definitively) unavailable when it does not react to the requests from the coordinator. The basic recovery policy for LH_{RS}^* is furthermore that (i) only the coordinator initiates a recovery, when a bucket does not react to the request from, and that (ii) the first action of any server upon its local start-up or restart is to report to the coordinator. The bucket then waits for the reply before providing any services, even if it finds its bucket intact.

Reliable communications and (i) prohibit any server from recovery while available. This excludes one possibility of the “bad” case. Next, the server restarting on its own could have its bucket already recovered by the coordinator, in the meantime, or not, if its unavailability remained unspotted. The coordinator systematically instructs any restarting server with the recovered bucket to become a spare. This precludes the bad case as well. If the unavailability remains unspotted, and the bucket is intact, the coordinator allows the bucket to serve. Otherwise, it finally initiates the recovery. The restarting server may eventually get the bucket back. Whatever the issue, the bad case cannot occur.

The scheme can be extended beyond the basic model, especially, to the case of unreliable communications that could temporarily or partly isolate a data bucket. The bucket could then become inaccessible to a client and the coordinator, while other clients could possibly continue accessing it. The coordinator would start the recovery that, if successful, could ultimately lead the clients using the former location to get stale data. To prevent this from happening, it suffices that any data server pings (scrubs) any of the parity buckets of its group at an interval shorter than a minimal bucket recovery time. In practice, as Section 5.7 shows, it should lead to a negligible overhead, for example, a ping per little bit more than a second in our experiments. If the server does not receive the reply, it alerts the coordinator, while possibly probing another parity bucket.

As will appear, every parity bucket knows the valid location of all the data buckets in its group, or knows about any recovery of a data bucket in progress in its group. In response to the ping, the available parity bucket should thus either confirm to the sender that it is still the available data bucket, or make it aware of the recovery in progress. Accordingly, the bucket should process any request coming afterwards, or during the ping, or should become a spare, only forwarding the misdirected requests to the coordinator. Notice that if the request is an insert, delete, or update, it must update all k parity buckets, as will be seen. These would redirect any such request to the valid bucket,

if a—necessarily unfinished—recovery had started since the latest ping. The misdirected key search, processed by the (invalid) data bucket alone, as any key search by any data bucket, cannot lead to an incorrect (stale) reply in this case either. No update could yet perform at the valid bucket, at least till the next ping.

2.2 Parity Records

The LH_{RS}* parity records enable the application to get the values of data records stored on any unavailable servers, up to $k \geq 1$. We call k the *availability level* and the file *k-available*. The actual level depends on the *intended availability level* $K \geq 1$ that is a file parameter. The K value adjusts dynamically to the file size, (Section 2.2.3). Typically, $k = K$, sometimes $k = K - 1$, after an increase of K . We now present how LH_{RS}* manages the parity records.

2.2.1 Record Grouping. LH_{RS}* parity records belong to the specific structure, invisible to the application. It consists of *bucket groups* and *record groups*. A bucket group contains m consecutive buckets with possibly fewer buckets in the last group. Formally, bucket group g consists of all buckets a such that $\lfloor a/m \rfloor = g$. Here, m is a file parameter that is a power of 2. Bucket group 0 consists of buckets $0 \dots m - 1$, bucket group 1 of buckets $m \dots 2m - 1$, and so forth. We limited our implementation to $m \leq 128$, since larger choices look to be of little use at present. Every data record in a bucket gets a unique *rank* $r = 1, 2, \dots$ when it enters the bucket because of an insert, split, or merge. Ranks are handed out basically successively, although the ranks of deleted records can get reused. The up to m data records sharing the same rank r in a bucket group g form a *record group* (g, r) . Each record group has $k \geq 1$ parity records stored respectively at a different *parity bucket* P_0, \dots, P_{k-1} . The value of k is the same for every record group in a bucket group and follows that of K as we have mentioned, and describe in depth in Section 2.2.3.

2.2.2 Record Structure. Figure 2b shows the structure of a parity record. Field r contains the rank and serves as the key. Field C encodes the record group structure. It contains m placeholders $c_0, c_1 \dots c_{m-1}$ for the keys of the data records in the group. If the i th-bucket in the group contains a data record with rank r and key c , then $c_i = c$ otherwise c_i is null. All k parity records in a record group share the value of r and of C . The final field is the *parity field* P , different for each parity record of a group. We generate it by encoding the D -fields in the record group using our Erasure Correcting Code (ECC) (Section 3). The ECC enables decoding any unavailable (hence assumed erased) $s \leq k$ D -fields in the group from any of $s \leq k$ parity records and the remaining $m - s$ D -fields. We can recover the unavailable s keys from the C -field of any parity record in the group. These properties are our basis for the k -availability.

In Litwin and Schwarz [2000], the actual keys for C formed a variable length list. The fixed structure above, proved more efficient [Ljungström 2000]. It typically needs slightly less storage. In addition, the position i of c in the C -field directly identifies the data bucket with c as the i th in its bucket group. This facilitates the search for data record c , that is needed with LH* addressing, for image management at the parity bucket and possibly for forwarding.

				77			
	69	2	47	[5,(--,77,--,--)]	[5,(--,77,--,--)]		
	37	7	7	[4,(--,69,--,47)]	[4,(--,69,--,47)]		
52	93	38	71	[3,(--,37, 2, 7)]	[3,(--,37, 2, 7)]		
4	61	22	87	[2,(52,93,38,71)]	[2,(52,93,38,71)]		
4	3	3	3	[1,(4,61,22,87)]	[1,(4,61,22,87)]		
				2,0	2,1		
4	5	6	7	P0	P1		

Fig. 3. LH_{RS}^* Record group structure with $k = 2$ parity buckets.

Example 3. We continue with our running example at Figure 1. We assume now that the bucket group size is $m = 4$. The file at Figure 1b then has four bucket groups $\{0, 1, 2, 3\}$, $\{4, 5, 6, 7\}$, $\{8, 9, 10, 11\}$, and $\{12\}$. Figure 3 shows the 2nd group that is group 1, with $k = 2$ parity buckets. These are named P0 and P1. We only show the keys in the data buckets, the ranks and the C -fields in the parity buckets. Notice that the latter are identical for both buckets, but that the P -fields would differ. The hyphens symbolize the null values. The header of a parity bucket only shows its bucket group number and the offset of the bucket in the group. The figure shows that the record group with rank 1 consists of m data records 4, 61, 22, and 87. The group of rank 3 contains only three records at present, in buckets 5, 6, 7. The next insert to bucket 3 will add a new member to this group, and will update all the C and P fields. If bucket 6 is unavailable, the data records in buckets 4, 5, 7 together with any of the parity buckets suffice to correct the erasure of the D fields in this bucket. The missing keys are in the C fields at the offset $i = 3$. To access record 93, among those necessary for the erasure correction calculus within its record group (1, 2), one may directly address bucket 5. Since $93 = c_1$ in its C -field, we calculate $1*4 + 1$. A 2-unavailability involving buckets 6, 7, requires the use of data buckets 4, 5 and of both parity buckets. A 3-unavailability is unrecoverable (catastrophic) here. We need at least $k = 3$ parity buckets to keep away from such bad luck.

2.2.3 Scalable Availability. The cost of storing and manipulating parity records increases with k . The storage overhead is at least k/m . We also need access to all the k parity records whenever we insert, update, or delete a data record. A multiple unavailability becomes more likely in a larger file, hence the probability of the catastrophic case for any given k rises as well, [Hellerstein et al. 1994]. In response, the LH_{RS}^* scheme provides *scalable availability* [Litwin et al. 1998]. At certain file sizes, the current availability levels increase by one for every group.

In more detail, we maintain a file parameter called the *intended availability level* K . Initially, $K = 1$. The coordinator increases K by 1, at the first split making the file extent N exceed some N_K data buckets, $K \geq 2$. Each N_K should be a power of 2 and a multiple of m . The former requirement makes K change only for the split of bucket 0. One choice for the latter requirement is $N_K = N_1^K$, with $N_1 = m$ [Litwin et al. 1998].

Consider now that the coordinator is at the point of increasing K since the file undergoes the split of bucket 0 making N exceed some N_K . At this moment, all bucket groups still contain $k = K - 1$ parity buckets. Now, whenever the

first bucket in a bucket group splits, its group gets one more parity bucket, the K -th one. From now on, every new split within the group updates this bucket as well. If the freshly appended bucket initiates a new group, then it gets K parity buckets from the start. Until the file reaches $N = 2N_K$ buckets, we thus have in the file some groups that are K -available while others are still $(K - 1)$ -available. The file is still $(K - 1)$ -available. Only when N reaches $2N_K$ will all the groups have K buckets, which makes the (entire) file K -available.

There is a subtle consequence for the *transitional* group that contains the current bucket n beyond its first bucket. Not all the buckets of the transitional group have yet split using h_{i+1} . Its K -th parity bucket therefore only encodes the data buckets in the group up to bucket $n-1$. LH_{RS}* recovery cannot use it in conjunction with the data buckets in the group that have not yet split. The group remains $(K - 1)$ -available despite K parity buckets. It reaches K -availability when the last bucket splits.

Example 4. We continue with $m = 4$. The file creation with $K = 1$ leads to data bucket 0 and one parity bucket for group 0, with buckets 0...3 as $m = 4$. The split creating bucket 4 initializes the first parity bucket for group 1, named P0 in Figure 3 and Figure 5, and formally named P_0 . Until the file reaches $N = 16$, every group has one parity bucket (P0), and the file is 1-available. The creation of bucket 16, by the split of bucket 0, appends P1 to group 0, and provides group 4 initiated by bucket 16, with P0 and P1 from the start. Group 0 remains transitional and 1-available, despite its two parity buckets. This lasts until bucket 3 splits. The new group 4 is however immediately 2-available.

Next, the eventual split of bucket 4 appends P2 to group 1, and starts group 5. Group 1 is transitional in turn. When bucket 15 splits, hence n returns to 0 and the file reaches $N = 32$ buckets, all the groups are 2-available. They remain so until the file reaches $N_3 = 64$ buckets. The next split, of bucket 0 necessarily, increases K to 3, adds P3 to group 0, makes the group transitional again, and so on. The file reaches 3-availability when it attains 128 data buckets. And so on.

3. PARITY CALCULUS

3.1 Overview

We first defined the parity calculus for LH_{RS}* in Litwin and Schwarz [2000] for our Erasure Correcting Code (ECC) derived from a classical Reed-Solomon code [MacWilliams and Sloane 1997]. We recall that an ECC encodes a vector \mathbf{a} of m data symbols into a *code word* of $m + k + k'$ code symbols such that any $m + k'$ of the code symbols suffices to recalculate \mathbf{a} . Our ECC was Maximum Distance Separable (MDS). Hence $k' = 0$ and m code symbols suffice to recalculate \mathbf{a} . This is the theoretical minimum, so we minimize parity storage overhead for any availability level k within a record group. Next, our ECC was systematic. Thus, our encoding concatenates a vector \mathbf{b} of k parity symbols to \mathbf{a} to form the code word $(\mathbf{a}|\mathbf{b})$. Finally, our code was linear. Thus, we calculate \mathbf{b} as $\mathbf{b} = \mathbf{a} \cdot \mathbf{P}$ with a *parity matrix* \mathbf{P} .

Our initial proposal was theoretical. Since then, we worked on the implementation. Our primary goal was fast processing of changes to data buckets.

The result was an improved ECC that we introduce now. We changed from symbols in the Galois field $\text{GF}(2^4)$ to $\text{GF}(2^8)$ or $\text{GF}(2^{16})$. We based our latest prototype [Litwin et al. 2004a, b], on the latter as measurements favored it. We also refined the parity matrix \mathbf{P} . Our new \mathbf{P} has a row and a column of ones. As a consequence of the latter, updating parity bucket P_0 involves only bit-wise XOR operations. In addition, the recovery of a single data bucket involves only the XOR operations as well. Our EEC calculus is thus in this case as simple and fast as that of the popular 1-available RAID schemes. The 1-unavailability is likely to be the most frequent case for our prospective applications. When we need to go beyond, the row of ones makes inserting, deleting, or updating the first data record in a record group, involve only XOR-operations at *all* parity records. Creation and maintenance of all other parity records involves symbol-wise GF-multiplication. It is mathematically impossible to find a parity matrix for an MDS code with more ones in it, so that our code is optimal in that regard. Our parity matrix is new for ECC theory, to the best of our knowledge.

We continue to implement our GF-multiplication using logarithms and antilogarithms [MacWilliams and Sloane 1997, Ch. 3, §4]. We take further advantage of this approach by directly applying to the parity (P -fields) generation, the logarithms of the entries in \mathbf{P} , instead of the entries themselves. The latter is the usual practice and we did so ourselves previously. The result speeds up the calculus, as we will show. Finally, we introduce the concept of a *generic parity matrix*. This matrix contains parity matrices \mathbf{P} for a variety of numbers m of data records, and k of parity records.

We now describe our present calculus in more depth. We only sketch aspects that our earlier publications already covered. See in particular Litwin et al. [2004b] for more details. We also give some details in Online Appendix C available in the ACM Digital Library.

3.2 Generic Parity Matrix

Our symbols are the 2^f bit strings of length f , $f = 8, 16$. We treat them as elements of $\text{GF}(2^f)$. The generic parity matrix \mathbf{P}_{gen} is the 2^{f-1} by $2^{f-1} + 1$ matrix of the form:

$$\mathbf{P}_{\text{gen}} = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ 1 & p_{1,1} & \cdots & p_{1,2^{f-1}} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & p_{2^{f-1}-1,1} & \cdots & p_{2^{f-1}-1,2^{f-1}} \end{pmatrix}. \quad (3.1)$$

With this choice of size, \mathbf{P}_{gen} can accommodate a maximum of 2^{f-1} data records and $2^{f-1} + 1$ parity records in a group. We could have chosen other maxima within the bounds of a total of $2^f + 1$ data and parity records in a group. We now sketch how to obtain \mathbf{P}_{gen} . We start with an extended Vandermonde matrix \mathbf{V} with $r = 2^{f-1}$ rows and $n = 2^f + 1$ columns over $\text{GF}(2^f)$. This is a matrix of largest known size so that any r by r sub-matrix is invertible. Elementary row transformations and multiplication of a column by a scalar preserve this property. We use them to transform \mathbf{V} into a matrix \mathbf{W} of the

form $(\mathbf{I} \mid *)$, a matrix whose left half is an identity matrix. We now multiply all rows j of \mathbf{W} with $w_{r,j}^{-1}$. Column r now only contains ones. However, the left r columns are no longer the identity matrix. Hence we multiply all columns $j \in \{0, \dots, r-1\}$ with $w_{r,j}$ to recoup the identity matrix in these columns. Next, we multiply all columns r, \dots, n with the inverse of the coefficient in the first row. The resulting matrix now also has 1-entries in the first row. This is our *generic* generator matrix \mathbf{G}_{gen} . Matrix \mathbf{P}_{gen} is its right half.

$$\mathbf{G}_{\text{gen}} = \begin{pmatrix} 1 & 0 & \cdots & 0 & 1 & 1 & \cdots & 1 \\ 0 & 1 & \cdots & 0 & 1 & p_{1,1} & \cdots & p_{1,r} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \cdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 1 & p_{r-1,1} & \cdots & p_{r-1,r} \end{pmatrix}. \quad (3.2)$$

Appendix A shows our actual generic parity matrices for GF(2⁸) and GF(2¹⁶). The Lemma in App. C in Online Appendix C available in the ACM Digital Library shows that any parity matrix obtained as the upper left corner of our generic parity matrix defines a systematic, linear, MDS ECC.

3.3 Parity Generation

We recall that the addition in a GF with 2^f elements is the bit-wise XOR operation. Also, among several methods for the multiplication, the popular method of logarithms and antilogarithms [MacWilliams and Sloane 1997] is especially convenient for our purpose. Given a primitive element $\alpha \in GF(2^f)$, we multiply two non-zero GF elements β and γ as $\beta \cdot \gamma = \text{antilog}_\alpha(\log_\alpha(\beta) + \log_\alpha(\gamma))$. We tabulate the logarithms and antilogarithms. There are 2^f - 1 entries in the logarithm table (one for each non-zero element) and twice as many in the antilogarithm table, one for each possible sum. By adding the logarithms modulo 2^f - 1, we could use a smaller antilogarithm table at the cost of speed. See App. C in Online Appendix C available in the ACM Digital Library for more details of our use of GF arithmetic, including the pseudo-codes. Table 5 there shows the logarithms based on $\alpha = 2$ we actually use for GF(2⁸).

We organize the nonkey fields of the m data records in a record group as the columns in an l by m matrix $\mathbf{A} = (a_{i,j})$, $0 \leq i < l$, $0 \leq j < m$. Similarly, we number the parity records in the group from 0 to $k - 1$ and arrange their P-fields as the columns in an l by k matrix $\mathbf{B} = (b_{i,j})$, $0 \leq i < l$, $0 \leq j < k$. Parity matrix \mathbf{P} is the upper left m by k corner of \mathbf{P}_{gen} . We define \mathbf{B} by $\mathbf{B} = \mathbf{A} \cdot \mathbf{P}$. Equivalently, we have:

$$b_{i,j} = \bigoplus_{v=0}^{m-1} a_{i,v} p_{v,j}. \quad (3.3)$$

Formula (3.3) defines the P -field of the parity records for a record group. We calculate this field operationally when we change a data record in the record group. This happens when an application inserts, deletes, or modifies a data record, or when a split or merge occurs. We implemented the latter as bulk inserts and deletes. For the parity calculus, inserts and deletions are special cases of updates. A missing data record has a nonkey field consisting of zeroes.

$$\mathbf{P} = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1a & 1c \\ 1 & 3b & 37 \\ 1 & ff & fd \end{pmatrix} \quad \mathbf{Q} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 105 & 200 \\ 0 & 120 & 185 \\ 0 & 175 & 80 \end{pmatrix}$$

Fig. 4. \mathbf{P} and \mathbf{Q} for 3-available record group of size four data record.

An insert changes this field from the zero string and a delete changes it to the zero string.

We now explain an update to data record j in a record group. Let $(\alpha_{i,j}^{\text{old}})_{0 \leq i < l}$ be the nonkey field of the record before, and $(\alpha_{i,j}^{\text{new}})_{0 \leq i < l}$ the one after, the update. The *delta field* (Δ -field) is $\Delta^j = (\delta_{i,j})_{0 \leq i < l} = (\alpha_{i,j}^{\text{new}} \oplus \alpha_{i,j}^{\text{old}})_{0 \leq i < l}$. It is the bit-wise XOR of the before and the after images of the nonkey field. For an insert or a delete, the Δ -field is thus simply the actual nonkey field. As in Litwin and Schwarz [2000], we calculate the \mathbf{P} -field of parity record s as:

$$b_{i,s}^{\text{new}} = b_{i,s}^{\text{old}} \oplus \delta_{i,s} \cdot p_{j,s}. \quad (3.4)$$

Parity bucket s only uses column s of \mathbf{P} . However, we do not use the parity matrix \mathbf{P} directly. Rather, we store the logarithms $q_{j,s} = \log_{\alpha}(p_{j,s})$ of the entries of \mathbf{P} in a logarithmic parity matrix \mathbf{Q} and use the formula

$$b_{i,s}^{\text{new}} = b_{i,s}^{\text{old}} \oplus \text{antilog}_{\alpha}(\log_{\alpha}(\delta_{i,s}) + q_{j,s}) \quad (3.5)$$

where the plus sign is the integer addition and the \oplus the bitwise XOR operation. Our use of \mathbf{Q} avoids half the look-ups to the logarithm table. Parity bucket s only needs the first m coefficients from column s of \mathbf{Q} , and only these are stored there. Parity bucket 0 does not need to store anything, since it simply updates according to

$$b_{i,s}^{\text{new}} = b_{i,s}^{\text{old}} \oplus \delta_{i,s}. \quad (3.6)$$

Appendix A shows our actual matrices \mathbf{Q} for $\text{GF}(2^8)$ and $\text{GF}(2^{16})$. At the implementation level, a data bucket calculates the Δ -field for each update and sends it together with the rank to each parity bucket. Each parity bucket then updates the \mathbf{P} -field of the parity record in the record group given by the rank according to (3.5), or to (3.6) if it is the first parity bucket.

Example 5. We use 1B symbols as elements of $\text{GF}(2^8)$ and write them as hexadecimal numbers. We continue with $m = 4$ and we choose $k = 3$. We cut the parity matrix \mathbf{P} and a logarithmic parity matrix \mathbf{Q} in Figure 4 from our generic parity matrix in Appendix A. We number the data buckets in a bucket group D0...D3 and consider three parity buckets, Figure 5. We have one data record per bucket. The record in D0 has the D -field: “En arche en o logos...”. The other D -fields are “In the beginning was the word...” in D1, “Au commencement était le mot...” in D2, and “Am Anfang war das Wort...” in D3. Assuming the ASCII coding, D -fields translate to (hex) GF symbols in Figure 5c, for example, “45 6e 20 61 72 63 68...” for the record in D0. We obtain the parity symbols in P0 from the vector $\mathbf{a}^0 = (45, 49, 41, 41)$ multiplied by \mathbf{P} . The result $\mathbf{b}^0 = \mathbf{a}^0 \cdot \mathbf{P}$ is (c, d2, d0). We calculate the first symbol of \mathbf{b}^0 simply as $45 \oplus 49 \oplus 41 \oplus 41 = c$.

D0	D1	D2	D3	P0	P1	P2
45	0	0	0	45	45	45
6e	0	0	0	6e	6e	6e
20	0	0	0	20	20	20
61	0	0	0	61	61	61
72	0	0	0	72	72	72
63	0	0	0	63	63	63
68	0	0	0	68	68	68

D0	D1	D2	D3	P0	P1	P2
45	49	0	0	c	41	ea
6e	6e	0	0	0	4b	32
20	20	0	0	0	47	87
61	70	0	0	11	75	48
72	72	0	0	0	52	63
63	69	0	0	a	0	6b
68	6e	0	0	6	4d	34

D0	D1	D2	D3	P0	P1	P2
45	49	41	44	9	f6	fe
6e	6e	6d	61	c	54	9
20	20	20	6e	4e	40	c1
61	70	41	73	23	d8	28
72	72	6e	20	4e	ce	4d
63	69	66	6c	0	18	39
68	6e	61	65	2	a0	a5

D0	D1	D2	D3	P0	P1	P2
49	49	41	44	5	fa	f2
6e	6e	6d	61	c	54	9
20	20	20	6e	4e	40	c1
74	70	41	73	36	cd	3d
68	72	6e	20	54	d4	57
65	69	66	6c	6	1e	3f
20	6e	61	65	4a	e8	ed

Fig. 5. Example of Parity Updating Calculus.

This is the conventional parity, as in a RAID. The second symbol of \mathbf{b}^0 is $45 \cdot 1 \oplus 49 \cdot 1a \oplus 41 \cdot 3b \oplus 41 \cdot ff$ and so on. Notice that we need no multiplication to calculate the first addend.

We now insert these records one by one. Figure 5a shows the result of inserting the first record into D0. The record is in fact replicated at all parity buckets, since updates to the first bucket translate to XOR operations at the parity buckets, and there are not yet the other data records. Inserting the second record into D1 at Figure 5b still leads to an XOR operation only at P0, but involves GF-multiplications using the respective \mathbf{P} columns at the other parity buckets. Notice that, in the latter case, we operationally use \mathbf{Q} columns only. Figure 5c shows the insert of the last two records. Finally, Figure 5d shows the result of changing the first record to “In the beginning was ...”. Operationally, we first calculate the Δ -field at D_0 : $(45 \oplus 49, 6e \oplus 6e, 20 \oplus 20, 61 \oplus 74, \dots) = (c, 0, 0, 15, \dots)$ and then forward it to all the parity buckets. Since this is an update to the first data bucket, we update the P-fields of *all* parity records by only XORing the Δ -field to the current contents.

3.4 Erasure Correction

Our ECC calculates a code word $(\mathbf{a} | \mathbf{b})$ as $(\mathbf{a} | \mathbf{b}) = \mathbf{a} \cdot \mathbf{G}$ with a *generator matrix* $\mathbf{G} = (\mathbf{I} | \mathbf{P})$ [MacWilliams and Sloane 1997, Ch. 1, §2], obtained by the concatenation of the m by m identity matrix \mathbf{I} and the parity matrix \mathbf{P} . We recall that we organized the non-key fields of the data records in an l by m matrix \mathbf{A} and similarly the P-fields of the parity records in a matrix $\mathbf{B} = \mathbf{A} \cdot \mathbf{P}$. Assume that we have m columns of $(\mathbf{A} | \mathbf{B})$ left, corresponding to m surviving records. We collect the corresponding columns of \mathbf{G} in an m by m submatrix \mathbf{H} . Here, the data buckets implicitly correspond to columns in the identity matrix. The columns of \mathbf{P} are reconstructed as the coordinate-wise antilogarithms of the columns in \mathbf{Q} at the parity buckets. We form an l by m matrix \mathbf{S} containing the m available data and parity records. We have $\mathbf{A} \cdot \mathbf{H} = \mathbf{S}$. Hence $\mathbf{A} = \mathbf{S} \cdot \mathbf{H}^{-1}$ and we recover all

$$\mathbf{H} = \begin{pmatrix} 0 & 1 & 1 & 1 \\ 0 & 1 & 1a & 1c \\ 0 & 1 & 3b & 37 \\ 1 & 1 & ff & fd \end{pmatrix} \quad \mathbf{H}^{-1} = \begin{pmatrix} 1 & a7 & a7 & 1 \\ 46 & 7a & 3d & 0 \\ 91 & c8 & 59 & 0 \\ d6 & b2 & 64 & 0 \end{pmatrix}$$

Fig. 6. Matrices for correcting erasure of buckets D0, D1, D2.

data records with one single matrix inversion (in our experiments, the classical Gaussian inversion appeared to be the fastest, and we actually use $\log_2(\mathbf{H}^{-1})$ to speed up the multiplication). If there are unavailable parity records, we recover the data records first, and then regenerate the erased P -fields by calculating $\mathbf{B} = \mathbf{A} \cdot \mathbf{P}$ for the lacking columns of \mathbf{B} .

Alternatively, we can decode all the erased data and parity values in a single pass. We form the *recovery matrix* $\mathbf{R} = \mathbf{H}^{-1} \cdot \mathbf{G}$. Since $\mathbf{S} = \mathbf{A} \cdot \mathbf{H}$, we have $\mathbf{A} = \mathbf{S} \cdot \mathbf{H}^{-1}$, hence $(\mathbf{A} | \mathbf{B}) = \mathbf{A} \cdot \mathbf{G} = \mathbf{S} \cdot \mathbf{H}^{-1} \cdot \mathbf{G} = \mathbf{S} \cdot \mathbf{R}$. We apply, however, only the first approach up to now. Computing \mathbf{R} is longer, while we expect primarily the recovery of the data bucket only.

Example 6. Assume that the data buckets D0, D1, and D2 in Figure 5 are unavailable. Assume that we want to read our record in D0. We collect the columns of \mathbf{G} corresponding to the available buckets D3, P0, P1, and P2 in Figure 5 in matrix \mathbf{H} , Figure 6. We invert \mathbf{H} . The last column of \mathbf{H}^{-1} is a unit vector since the fourth data record is among the available ones. To reconstruct the first symbols simultaneously in each data bucket, we form a vector \mathbf{s} from the first symbols in the available records of D3, P0, P1, and P2: $\mathbf{s} = (44, 5, fa, f2)$. This vector is the first row of \mathbf{S} . To recover the first symbol in D0, we multiply \mathbf{s} by the first column of \mathbf{H}^{-1} and obtain $49 = 1 \cdot 44 + 46 \cdot 5 + 91 \cdot fa + d6 \cdot f2$. Notice again that we actually use the matrix $\log_2(\mathbf{H}^{-1})$. We iterate over the other rows of \mathbf{S} to obtain the other symbols in D0. If we were to read our record in D1, we would use \mathbf{S} with the second column of \mathbf{H}^{-1} .

4. LH_{RS}^* FILE MANIPULATION

An application manipulates an LH_{RS}^* file as an LH^* file. The coordinator manages high-availability invisibly to the application. We assume that the coordinator can diagnose any bucket unavailability by *probing* through reliable communications. Available buckets respond to the probing as expected, by definition. Internally, each manipulation starts in normal mode. It remains so as long as it does not run into an unavailable bucket. That one does not respond at all (our basic case), or more generally does not behave as it should. A manipulation that encounters an unavailable bucket enters *degraded mode*. The node handling the manipulation forwards it to the coordinator. The coordinator initiates the *bucket recovery*, unless it is already in progress. It considers the bucket content erased and calls upon our erasure correction calculus. The recovery simultaneously restores *all* unavailable buckets found in the group when the operation starts, up to k , as we discussed. More than k unavailable buckets make the unavailability catastrophic.

For a key search of a record in an unavailable bucket, the coordinator also starts the *record recovery*. This operation recovers only the record, or finds its

key not to be in the file. It should usually be much faster than waiting for the bucket recovery to complete, which also recovers the record. Once the bucket recovery, or only the record recovery for the key search, has terminated, the coordinator completes the original manipulation and returns the result and the control to the requestor.

We first present the recovery operations. Next, we describe the record insert, update and bucket split operations. We focus on the parity coding, as Section 2 already largely covered the related data record manipulations. We have left the description of file creation and removal, of key search, and of nonkey search (scan) to Online Appendix C available in the ACM Digital Library, App. C. These manipulations do not involve parity coding and perform almost as for LH*, except for the degraded mode. Likewise, we relegate the deletion and bucket merge to Online Appendix C available in the ACM Digital Library. Deletions are rare (except those internal to the split operation we describe in Section 4.5), while merges are hardly ever implemented.

4.1 Bucket Recovery

The coordinator starts the bucket recovery by probing the m data and k parity buckets of the bucket group for availability. We usually have $k = K$, or we have $k = K - 1$ if the group is not aware yet of the current K , or is the transitional one. The probing should localize m available buckets, excluding the last parity bucket in a transitional group. The unavailability is catastrophic otherwise. The coordinator reports the gloomy fact to the client, for the application, and stops the processing. It may still be possible to recover some records, but this is beyond the basic scheme.

The probing also localizes the unavailable buckets that should be therefore k at most, including the originally reported bucket, and without the last bucket of the transitional group.⁴ If $k > 1$, then prior to the erasure correction, the coordinator may need to verify whether all the available parity buckets processed the last Δ -record or even simply received it. It may not be the case, as we show in Section 4.3 and Section 4.4. If Δ -record could update some, but not all, parity buckets, the erasure correction result could be gibberish, corrupting the data. The coordinator first terminates the updates at those parity buckets that have updates pending. Next, the coordinator establishes the list L_A of m available buckets. L_A possibly includes bucket P_0 . The coordinator establishes also the list L_S of tuples: (unavailable bucket, replacing spare). The coordinator chooses one spare as the *recovery manager*. It passes the task to it. If there is any parity bucket in L_S , it also passes its \mathbf{Q} column. The handover prevents the coordinator from becoming a hot spot.

The manager first creates the empty, structure of the bucket to recover at each spare. Next, if there is only one unavailable (data) bucket, and $P_0 \in L_A$, then the manager applies the XOR-only erasure correction. Otherwise, it creates matrix \mathbf{H} using the columns of \mathbf{Q} at the parity buckets in L_A . If the bucket

⁴In the basic scheme, we exclude the case of the probing finding a bucket reported to be unavailable to be nevertheless available, because for example only the reporting client had problems communicating with the bucket.

group is the last one, and some data buckets have not yet been created, then the calculus considers these buckets as virtual ones with records having zero D -fields. In both cases, the manager then calculates the matrix $\log_2(\mathbf{H}^{-1})$, using the primitive element $\alpha = 2$ and the Gaussian inversion to produce \mathbf{H}^{-1} , as we said in Section 3.4. Next, it starts looping over all the parity records in a parity bucket in L_A . It requests the successive records, and for each record received, it performs the *record group* recovery, producing all the unavailable records of one group.

Record group recovery explores the C -fields, Figure 2a. For every key $c_i \in C$, it requests the data record c_i from the i th bucket in the group, provided it is in L_A . The manager decodes the (erased) D -fields of unavailable data records in the group. It uses the XOR-only or $\log_2(\mathbf{H}^{-1})$. It also reconstructs the keys from C . If there are also unavailable parity buckets in L_S , then the manager generates their records from the m data records. The procedure varies slightly for the last parity bucket in a transitional group. The recovery of this bucket uses only the data buckets among the m up to bucket n (nonincluded). Finally, the manager sends the recovered records to the spares.

Once the bucket group recovery ends, the manager sends the addresses of the recovered buckets to the remaining buckets in the group. These update the location tables. The manager finally returns control to the coordinator. The coordinator updates its server addresses as well. Clients and servers get the new addresses when they misdirect the requests.

It is perhaps worth recalling furthermore that the erasure correction used for each record group recovery could reconstruct the data and parity buckets simultaneously, except for the last bucket of the transitional group. It could use the alternate erasure correction calculus discussed in Section 3.4. The unavailability involving only the data buckets may be expected however, to be more frequent than that of both types together. Our current erasure correction is then faster.

While the bucket group recovery loops over the individual record recovery, the actual records sent among the data buckets, the parity buckets, and the spare, move in bulk. At present, we use the TCP/IP in passive mode, Section 5.1. This turns out to be far more effective than our previous choice of UDP with a record per datagram.

4.2 Record Recovery

A record recovery results from a key search in an unavailable correct bucket a . It decodes only the requested record or finds that the key is not in the file. This suffices for completing the search and should be typically much faster than delaying the key search until the sufficient completion of the bucket recovery, (Section 5.7 and Section 5.8). The coordinator starts the record recovery in parallel to a bucket recovery, provided the unavailability is recoverable, of course. It hands control to the *record recovery manager* at an available parity bucket in the group. It transmits c, a, L_A . The manager first searches the C -fields in the bucket for the existence of a parity record with rank r and the searched key c in $C(r)$. Since one cannot infer r from c , it basically scans the bucket. It visits in each record only the C -field value at the offset of the unavailable bucket in

the group. An index at the parity bucket, giving the ranks of the keys in the C -fields, could obviously avoid the scan, at the cost of additional storage and processing. Such an index is not part of the basic scheme.

If the manager does not find the parity record, it informs the coordinator that the search is unsuccessful. Otherwise, the recovery manager uses the erasure correction on the record group r to recover the searched D -field. As already discussed, the calculus uses possibly XOR only, or \mathbf{H}^{-1} otherwise. Unlike for the record group recovery within the bucket recovery, the record recovery manager restores only the requested record, even if there are more unavailable data records in the group. Finally, it sends the record to the coordinator that forwards it to the client.

4.3 Insert

In normal mode, an LH_{RS}* client performs an insert like an LH* client. The client then addresses the insert to the bucket determined by the data record key c and the client's image. It keeps the copy of the record and waits for an acknowledgement. If it does not come within a timeout, the client sends the insert to the coordinator. The operation enters the degraded mode. We have already presented the bucket recovery that it starts. Later in this section we address the aspects of the degraded mode specific to the insert. The client waits for the final acknowledgment to discard its copy of the record. The policy for its acknowledgements to the application is implementation dependent and is not within our basic scheme. Our prototype uses its flow control algorithm.

The bucket receiving the request eventually forwards it as described before. Once the correct data bucket receives the insert, it stores the record as for an LH* file. If the data bucket overflows, the bucket informs the coordinator. In addition, it assigns a rank r to the record. Next, it sends the Δ -record (with key) c and r to the k parity buckets. Recall that the Δ -field is the D -field of the inserted record. The data bucket then commits (acknowledges) the insert to the client and waits, internally, for the k acknowledgements from all parity buckets. The client discards its copy of the record.

The policy provides the k -availability under reliable messaging and our other basic assumptions, provided that all the available parity buckets perform the required update, as detailed below. The data bucket will report any unavailability and will still have the Δ -record. If the probing during the recovery finds any $l \leq k$ parity buckets available, all l will be updated with the Δ -record. If the probing also shows at most l data buckets unavailable, then our basic erasure correction calculus first recovers all of them and next the unavailable parity buckets. This, even if the probing finds even the original data bucket among the unavailable ones, as it could become unavailable in the meantime.

To update its content upon the message from the data bucket, each parity bucket creates the parity record r , if it does not exist already, and inserts c into C -field. It also encodes the Δ -field of Δ -record c through the update of P -field of record r . It first conceptually multiplies the Δ -field symbol-wise with the correct coefficient of matrix \mathbf{P} . The actual calculus does not use the multiplication for the first column as it has 1's only, and uses the symbols in \mathbf{Q} for $k > 1$. It either

XORs the result to the P -field already in record r , or stores it as the new P -field of new record r .

If $k = 1$, the send-out of the Δ -record amounts in our scheme to 1PC. That is, the parity bucket creates or updates the parity record r , acknowledges the operation, and does not keep any other traces of it. The data bucket also erases any traces after the acknowledgement to the client. That one does the same, perhaps after forwarding the acknowledgement to the application in turn. The server or the client enters the degraded mode when any of the expected messages does not arrive in time.

The approach does not suffice in contrast to provide k -availability for $k > 1$. For instance, assume that $k = 2$, and that both the data bucket and the client become unavailable, while the data bucket is sending the messages to the parity buckets. The unavailability of the data bucket can then be discovered only later on. It may then happen that the bucket could send the message only to one of the parity buckets that performed the creation/update of record r consequently. Consider that the probing finds yet another data bucket unavailable. The erasure correction calculus obviously cannot recover the group anymore, leading to data corruption in the case of such an attempt. This includes the inserted record itself that, besides, is nowhere around anymore. In any case, the group, hence the file, is not k -available anymore, despite the k parity buckets in the group. Notice that for an update, the case could occur even if the data bucket only became unavailable during the send-out. Unlike for an insert or delete, the C -field of a parity bucket does not suffice anymore to know whether it dealt with the Δ -record or not.

The general rule for $k > 1$ that follows from this, is that a change, whether insert, update or delete, should commit at all or none of the parity buckets. Our basic scheme uses the following simple variant of 2PC for any of these operations. The data bucket sends the Δ -record c and its r to all k parity buckets. It does it in the order of the column index in \mathbf{P} : $p_0, p_1 \dots p_k$. Each parity bucket starts the commit process by acknowledging the reception of the message. The confirmation constitutes the “ready-to-commit” message. Each parity bucket encodes the record as usual, into parity record r . But it retains the Δ -record in a differential file (buffer) for a possible rollback. If the data bucket gets all k “ready-to-commit” messages, it sends out the “commit” message to the k buckets, in the same order. Each bucket that receives the message discards the Δ -record. All available buckets should receive it under our general assumptions.

Degraded mode starts when any of the buckets involved cannot get a response it expects. The operation at the data bucket enters degraded mode if it lacks any of the acknowledgments from the parity buckets. It alerts the coordinator, transmitting the Δ -record, r , and as usual the number p , of the unavailable parity bucket. Within the bucket recovery process as described above, the coordinator requests each available bucket to complete the update from the differential file, or by using the Δ -record it sends out otherwise. The recovery process then proceeds further as already presented.

Another degraded case occurs when parity bucket p_l does not receive a message from the data bucket. The data bucket must have then just failed and bucket p_l must be in the “Ready-to-Commit” state and must still have the

Δ -record. It alerts the coordinator, sending out the Δ -record, and r . The coordinator probes the parity buckets in the above order. Any bucket up to p_l either has committed or still has the Δ -record. Any other bucket either still has the Δ -record or did not get it at all. The coordinator sends the Δ -record where needed and requests the update at all the available parity buckets accordingly. Again, the recovery process continues as already presented.

Next, the client lacking the acknowledgement from the data bucket reports it as unavailable. The client also acknowledges the operation to the application, to avoid further waiting. The coordinator performs the delivery to the correct bucket if the unavailable one wasn't. This may generate a split processed as usual. The coordinator may also find that the correct bucket it found on its own is unavailable as well. This may trigger a separate recovery if the bucket is in a different bucket group than the one reported by the client. Next, the coordinator determines the availability of buckets in the group, and verifies the synchronization of the parity buckets, as just described. It may find that the data bucket never sent any messages to parity buckets. Alternatively, it may also find an alert from a parity bucket about the data bucket's unavailability. It may further find some parity buckets still waiting for the commit from the data bucket and some perhaps without any message. In all cases, the coordinator can determine the state of each available parity bucket and synchronize all of them as above. Afterwards the recovery proceeds as described.

Finally, as we said, it may happen that the client and the data bucket become simultaneously unavailable during the insert. If the update was in progress, an available parity bucket would alert the coordinator as described. Alternatively, it may only be that no (available) parity bucket has received the Δ -record or all have nicely committed. Only a later application operation will discover the data bucket unavailability. It will be accordingly recovered with or without the inserted record.

An insert in a degraded mode to the unavailable correct data bucket may generate an overflow at the recovered bucket. The new bucket itself alerts the coordinator to perform a split. On the other hand, observe that one can optimize the 2PC to use multicast messaging to the parity buckets for inserts and deletes. It is due to the possibility of recovery synchronization using the C -fields.

Finally, notice that, while our 2 PC is a sure solution for k -availability and 1 PC cannot be, the event of a data bucket becoming unavailable during the messaging to the parity buckets is very unlikely. Section 5.4 shows this time to be under a millisecond. The bad case of both the client and the data bucket becoming unavailable during that time is obviously even much more remote. An application may be reasonably tempted to use 1PC anyhow, for its simplicity and necessarily better performance. After all, every application already disregards scores of potential, but fortunately very unlikely, causes of data unavailability.

4.4 Update

An update operation of record c changes its nonkey field. In the normal mode, the client performs the update as in LH*. As for an insert, it waits for an acknowledgement. The client sends the received record with its key c and the new

value of the D -field. The data bucket uses c to look up the record, determines its rank r , calculates the Δ -record, and sends both to all k parity buckets. These recalculate the parity records. Finally, the data bucket commits the operation.

As for inserts and deletes, 1PC suffices only for $k = 1$. To be on the safe side for $k > 1$, the updates also basically use the 2PC above. It clearly suffices for this operation as well.

4.5 Split

As in LH^* , if an insert to an LH_{RS}^* data bucket a results in an overflow, then the bucket a notifies the coordinator. The coordinator starts the *split* operation of bucket n , determined by the split pointer. We recall that this operation is invisible to the application. Typically, we have $n \neq a$. In the normal mode, the coordinator first locates an available server and allocates there the new data bucket N , where N denotes the number of data buckets in the file before the split. Bucket N is usually in a bucket group different from that of bucket n , unless the file is small and $N < m$. If N is the first bucket in the group, then the coordinator allocates K empty parity buckets. If $K > k$ for the bucket group with bucket n , then the coordinator allocates the additional K th parity bucket. Provided all this performs normally, the coordinator sends the *split* message to bucket n with all the addresses. This hands control of the split to bucket n . The coordinator nevertheless waits for the commit message. The bucket sends all the data records that change address when rehashed using h_{j+1} to data bucket N . Our implementation sends these records in bulk.

For each data record that moves, bucket n finds its rank r , produces a Δ -record that is actually identical to the record itself, and requests its deletion from the parity records r in all the k buckets of its group. It also assigns new successive ranks r' starting from $r' = 1$ to the remaining data records. Bucket n then sends both ranks with each Δ -record to the K parity buckets. Each existing bucket, deletes Δ -record from parity record r and inserts into parity record r' . The new K th parity bucket, if there is one, disregards the deletes.

When data bucket N receives the data records, it requests the insert into its K parity buckets with the successive ranks it assigns. Once it terminates, it reports to bucket n that in turn reports to the coordinator.

The operations on the parity buckets use 1PC for $K = 1$ and the already presented 2PC otherwise. Degraded mode starts when a data or a parity bucket does not reply. We skip the discussion of this mode here, as various cases are similar to those already discussed. Once the split terminates, the coordinator adjusts the file state as in Section 2.1.1.

5. PERFORMANCE ANALYSIS

We have analyzed storage, communication, and processing performance of the scheme. We have derived formulae for the load factor, parity storage overhead, and messaging costs. We limited the derivation to the dominant cost factors. This makes the calculus easy enough, but still quite lengthy. This analysis is in Online Appendix C available in the ACM Digital Library App. C.

The results show that the high availability of an LH_{RS}^* file incurs about the smallest possible storage overhead. More precisely, for any intended availability

level K , and of group size m , the load factor of the growing LH_{RS}* file should be in practice about constant, and the highest possible for these values, as well as for any technique added to an LH* file to make it K -available. Through motivating examples in Online Appendix C available in the ACM Digital Library, we show some practical design choices.

A user is mainly interested in the response time of various operations. The complexity of any practical implementation of LH_{RS}* seems to prevent a practically useful formal analysis of such times. We have preferred the experimental analysis of various implementations. We now present the results. They complete the picture of the efficiency of our scheme and they validate the various design choices that we presented above.

5.1 Prototyping LH_{RS}*

We have implemented LH_{RS}* to measure various operations and prove the viability of the scheme. The work took many years of effort. The earliest prototype is presented in Ljungström [2000]. It implemented the parity calculus defined in Litwin and Schwarz [2000]. It also reused an LH_{LH}* implementation for the data bucket management [Bennour 2002]. Experiments with the next version of the LH_{RS}* prototype are in Moussa and Litwin [2002]. The current version used for the experiments below builds upon that one. We present the prototype in greater detail in Litwin et al. [2004a, b], as we have shown it to the public at VLDB-04. Further details, as well as the deeper discussion of the experiments discussed below, are in Moussa and Litwin [2002].

The prototype consists of the LH_{RS}* client and server nodes. These are C++ programs running under the Windows 2000 Server. Internally, each client and server processes queries and data using threads. The threads communicate through queues and other data structures and synchronize on events. We mainly use two kinds of threads. The *listening threads* manage the communications at each node. There is one thread for UDP, one for TCP/IP and one for multicast messaging. The *working threads* (currently four) simultaneously process queries and data, whether received or sent out.

The communication uses the standard UDP and TCP/IP protocols. Clients communicate with servers through UDP, except when the data records are larger than a datagram could be (64 KB). A listening thread timely unloads the UDP buffers to prevent losing a datagram. There is also a flow control mechanism. The servers communicate using TCP/IP for data transmission during the bucket split or recovery, and UDP for other purposes. Again, a listening thread unloads the UDP buffers. Another such thread manages the TCP/IP stack. This stack has the listening socket in passive open mode [www.faqs.org/rfcs/rfc793.html]. This new connection mode, available in Windows 2000 [Litwin et al. 2004b], replaced those studied in our earlier prototypes. It handles a larger number of incoming requests more effectively. In fact, it skips the connection dialogue that previously was necessary for each request. It proved to be by far the most efficient connection mode in our experience with LH_{RS}*.

Many measures of the operations using the parity calculus reported below compare the use of $GF(2^8)$ and of $GF(2^{16})$. We expected the latter to be faster.

Experiments mostly (but not always) confirmed our intuition and quantified it. $GF(2^{16})$ provided the most noticeable acceleration for the erasure correction. We could also confirm and measure the benefit of using our logarithmic matrix \mathbf{Q} , derived from our newest version of a parity matrix \mathbf{P} , with a first column and a first row of ones. We then measured the speed of the operations involving the parity updates, namely the inserts, file creation with splits, and updates, as well as bucket and record recovery. The study varied the availability level from 0 to 3. We left the study of delete, merge and scan operations for the future. As we already said, the first two operations are of lesser practical interest, whereas the latter is independent of the parity calculus unless it triggers a record recovery. We also measured the speed of key searches as a basic reference. We averaged each measure over several independent experiments.

Practical considerations lead to simplified implementation of some operations. Also, the experiments modified our own ideas on the best design of some operations. We discuss these issues in the respective sections.

The test-bed for our experiments included five P4 PCs with 1.8 GHz clock rate and 512 MB memory, and a 2.6 GHz, 512 MB P4 machine. We used the latter as a client. Others served as data and parity servers. Sometimes, we also used additional client machines (733 MHz, P3). Our network was a 1 Gbps Ethernet.

5.2 Parity Generation

To test the efficacy of using \mathbf{Q} , we conducted experiments creating parity records in a bucket with a logarithmic \mathbf{Q} column, versus its original \mathbf{P} column. We used a group of $m = 4$ data buckets and created a parity bucket using the second or third or fourth parity column of each matrix (the first column of \mathbf{P} was that of ones). A data bucket contained 31250 records. Using $GF(2^8)$, the average processing time shrank from 1.809 sec to 1.721 sec. We saved 4.86%. Use of $GF(2^{16})$, reduced the time from 1.462 sec to 1.412 sec, that is, by 3.42%. Notice that $GF(2^{16})$ was always faster, by about 20%.

We investigated the influence of the column and row of ones in the parity matrix \mathbf{P} , or equivalently, of a column and row of zeroes in \mathbf{Q} . This means that updates to the first data bucket only involve XORing, but no Galois field multiplications. For $GF(2^8)$, the processing time shrank further from 1.721 sec to 1.606 sec, that is, by 6.68%. Using $GF(2^{16})$, we measured 1.412 sec and 1.359 sec, that is, 3.75% of additional savings. Again, $GF(2^{16})$ was always faster, but by only about 15%.

As expected the \mathbf{Q} with first column and first row of zeroes yields the fastest encoding. We therefore used only this choice for our experiments below. We attribute the always higher savings for $GF(2^8)$ to the higher efficacy of XORing byte-sized symbols.

5.3 Key Search

The key search times in the normal mode serve as the reference for the access performance of the prototype, since they do not involve the parity calculus. We measured the time to perform random *individual* (synchronous) and *bulk*

(asynchronous) successful key searches. All measurements were done at the client. We start the timing of an individual search when the client gets the key from the application. It ends when the client returns the record received from the correct server. The search time reported is the average over a synchronous series of individual searches, one starting after the end of another. We start the clock to measure a bulk search when the client gets the first key from the application and we stop it when the application receives the last record searched. We report the average time. During the bulk search, the client launches searches asynchronously, as usual for a database query. These searches use UDP and a custom flow control method to prevent a server overload. Interestingly, most of our experiments did not even trigger this mechanism.

We measured the search time in a file of 125,000 records, distributed over four buckets and servers. A record had a 4 B key and 100 B of nonkey data. The average individual search time was 0.24 ms. The bulk one was 0.06 ms, four times faster. The individual search is thus about 40 times faster than a single disk access, whereas the bulk one is about 200 times faster. The server CPU speed was the limiting factor for the former and the speed of the client for the later.

5.4 Insert

We timed a series of individual and bulk inserts in the normal mode. The inserts addressed bucket 0, without triggering any bucket split. The idea was to test the most unfavorable scenario, in which all the inserts from a client end up at a single bucket. Inserts into multiple buckets with splits gave rise to different experiments, Section 5.5. We defined the individual insert time for this experiment as starting when the client receives the record from the application and ending when the client synchronously gets the acknowledgement from the data bucket. Bulk inserts used asynchronous acknowledgements for the flow control, as for the bulk key searches. The client acknowledges the insert to the application after it gets its own acknowledgement. Otherwise an overflow of its queues could result.

We did not measure the degraded mode. It would amount to collecting arbitrary timeouts. We only implemented 1PC, leaving the more complex 2PC protocol for the future. As we mentioned already, some, perhaps many, prospective implementations are also likely to choose only 1PC. We have nevertheless made additional measurements, forecasting the 2PC performance as well. Our experiments for the inserts measured the basic case of the data bucket sending the acknowledgement to the client immediately after the messages to its k parity buckets. Using 2PC would not change the insert time under this condition, as long as the additional load did not saturate the server.

We timed a series of 10,000 inserts into an initially empty bucket of $b = 10,000$, thus avoiding the split, as wished. Again, a record consisted of a 4 B key and 100 B of nonkey data. We recorded 0.29 ms for $k = 0$, 0.33 ms for $k = 1$ and 0.36 ms for $k = 2$. The choice of GF did not matter as the Δ -record is simply the inserted one. The average bulk insert time was 0.04 ms, seven to nine times faster. That time was the same as for updates, as discussed in Section 5.6. Both

Table I. Average Bulk and Individual Blind Update Times (Milliseconds) per Record

	Bulk	Individual			
		$k = 0$	$k = 1$	$k = 2$	$k = 3$
$GF(2^8)$	0.04	0.25	0.48	0.57	0.58
$GF(2^{16})$	0.04	0.24	0.50	0.55	0.59

bulk times were measured in the same way, and are similarly independent of k .

The previous figures show that adding the first parity bucket to a 0-available file slows down an insert on the average by 0.04 ms, or 14%. Adding one more parity bucket costs slightly less, 0.03 ms, or 10%. The increment is due mostly to the additional message. Its cost is about $0.03 \div 0.04$ ms, as the bulk insert time shows as well. This allowed us to trivially infer the timing for larger k 's.

Table I confirms these calculations for the updates. It also shows that the time for the message with the Δ -record followed by its acknowledgment was, in those experiments, about 0.05, about the double of the time above. These numbers allow us to estimate the total time of an insert-processing at the data bucket, and the throughput, if 2PC was used. Related details are in Section 5.6. Applied here, for instance to $k = 3$, they show this time to be about $0.36 + 4 \cdot 0.03 = 0.48$ ms. The rationale for the formula is the time for the last acknowledgment yet to come in from the server (others came normally earlier, in parallel to the outgoing messages from the client), followed by three commit messages sent in order by the client. Accordingly, we forecast the throughput as about 2000 inserts/s. For $k = 2$ the time would be smaller, about 0.42 ms, and the throughput higher, reaching about 2400 inserts/s. And so on.

All this appears to be a quite efficient behavior. Notice finally that the measured insert times at the client are respectively at least about 30 to 250 times faster than to local disks (assuming 10 ms per access). As for a key search, the individual insert time was bound mainly by the server speed, while the bulk insert time was bound by the maximal client speed.

5.5 File Creation

Figure 7 shows the average file creation time, by inserts with splits this time, for a bucket group of $m = 4$ data buckets and $k = 0, 1, 2$ parity buckets. The inserts are individual ones. We did not experiment with bulk inserts, as they need a more complex design of splits, left for future work, to prevent side effects resulting from the concurrent processing of splits and of inserts. Besides, the average time to create a file using l record bulk inserts would be at the client, simply $0.04l$ ms, given the bulk insert time. At a server, the time could be longer, to complete the last inserts (see the following discussion of bulk updates). For the previous experiments with inserts, during the file creation the data bucket sends the acknowledgement to the client, after sending the messages to the k parity buckets, but without waiting for the acknowledgements from these buckets. The results we measured were practically the same for $GF(2^{16})$ and $GF(2^8)$. Hence, the charts shown apply to both fields, although we give the

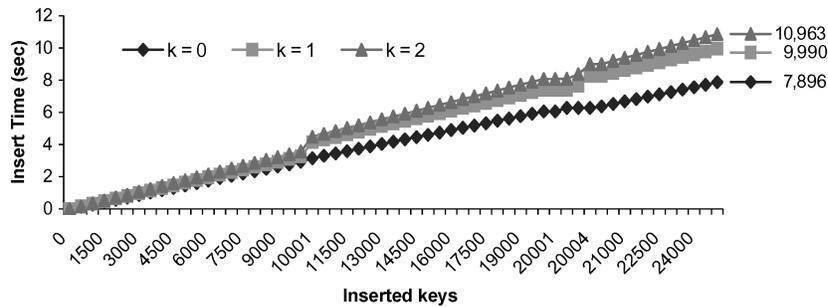


Fig. 7. Creation times (seconds).

numerical values for $GF(2^8)$. We inserted a series of 25,000 records, again with a 4 B key and 100 B of nonkey data per record. The bucket size was $b = 10,000$. A point of the chart corresponding to l inserts shows the total time to perform these inserts.

The inserts caused the file to split thrice. The split of bucket 0 occurred naturally after insert 10,000. A temporary slowdown of the insert times resulted, greater for greater k . The next inserts went uniformly into buckets 0 and 1. After slightly more than 10,000 further inserts, both buckets split almost concurrently. That is why the chart seems to show only two splits.

From the times to insert the 25,000 records, in the figure, we can gauge the typical cost of additional parity buckets for our file, once it scales to the steady state with many groups. For $k = 0$, we have the creation time of 7.985 s. For $k = 1$, we have 10.125 s, that is 27% more. Finally, for $k = 2$, the time is 10.974 s, that is 8% slower than for $k = 1$. The related average times per record inserted were 0.32 ms, 0.41 ms, and 0.44 ms for $k = 0, 1, 2$ respectively. Splits thus respectively introduced additional average costs of 3 and 8 ms, as compared to the costs of individual inserts alone. The percentage values are respectively 10.4%, and 24%. All together, these times are at least 20–30 times faster than disk accesses.

As is to be expected, adding the first parity bucket causes the most noticeable degradation. The percentage value of 27% is about twice that for an insert alone. We now see a cost of the parity calculus for the splits. Adding parity buckets has globally much less effect (a 8% slow-down), also because of the parallelism of the parity updates. Notice however that there is no incidence on the cost of the updates to the new parity bucket during the splits, as the difference to the average time per insert for $k = 2$ remains 8 ms. It confirms logically that split processing on the parity buckets is about fully parallel. We extrapolate the increase for each value of $k > 2$ to be the same 8%. The increase is caused mainly by the additional messaging at the data bucket.

The charts in the figure are about linear. The experiments thus confirm the scalability of the scheme, and we can predict the creation times for larger files. We create our files for $k = 0, 1, 2$ at the rate (speed), respectively, of 3131, 2469 and of 2278 records per second. For instance, to scale up our 2-available file to 1 M records should thus take 439 seconds about 7.3 minutes. More generally, as our records are 104 B long, we create our files at a rate of 0.33 MB/sec for $k = 0$,

0.25 MB/sec for $k = 1$, and of 0.23 MB/sec for $k = 2$. These numbers allow us to predict linear creation times for other record sizes. Bulk creation times and rates should be limited by the client and yet be about ten times faster from the application's point of view. For instance, less than a minute should suffice for a 1 M record file.

We also timed the use of our former \mathbf{Q} matrix, without the first column and row of ones. The creation time for $k = 1$ was 10.011 sec. Thus, our new \mathbf{Q} speeds up encoding time by almost 2%. While the acceleration appears to be slight, we recall that it comes at no cost.

5.6 Update

To determine the update performance, we generated series of 500, 1000, 5000, and 8000 blind updates to the records in our LH_{RS}^* file (same as for the insert experiments). We updated different records, to prevent the caching. Table I shows the results for bulk and individual updates. All updates used UDP and 1PC. As before, we only measured 1PC. This time, however, the data bucket waits for all the acknowledgements before sending the commitment to the client. The second column gives the average bulk update time in the normal mode. These measures start with the reception of the first update from the application and ends with the send-out of the last of the series. The processing at the servers may last longer. In addition, if the series is longer, then more records are perhaps temporarily stored in the queue of the listening thread at each server. Some acknowledgements may come back to the client after the end of the bulk update. The processing time at the data bucket depends on k . Nevertheless, it does not influence the bulk update time as defined here. If any acknowledgements were negative, or missing, the client would start the degraded mode. Notice that the bulk insert time is independent of the GF used.

The other columns list the average individual update times for $k = 0 \dots 3$. The bulk update times are basically six times faster, as the comparison for $k = 0$ shows. The numbers show also that using one parity bucket doubles the update time into the data bucket alone. This result matches the intuition. However, adding more parity buckets only increases the time by 10% to 20%. Notice that adding the 3rd parity bucket adds only 2–7%. All this is again, nice behavior. One may further extrapolate these results to the server processing of the bulk updates. Finally, using $GF(2^{16})$ does not appear uniformly faster. The results are practically identical for both fields, as for inserts.

Compared to the inserts, the bulk times also do not change, as the client processes inserts and updates at the same possible speed. The individual update time takes, in contrast, considerably longer. The times in Table I for $k = 1$ are already almost 45% longer than the time to insert. It is the measure of the additional processing of the Δ -record (XORing with the before image) followed by the UDP messaging and the waiting for the acknowledgements. The first ones come back in parallel to the outgoing messages, hence only the last one actually counts. To speed up the parallelism, the listening thread at the parity buckets acknowledges the Δ -record at the earliest moment, namely when it saves it in its queue from the communication buffer, before the actual encoding into the

parity record. That is why the times to process the message from the bucket followed by its acknowledgment are only about twice of the cost of a message, about 0.05 ms on the average. This matches the results for the inserts where the measured message time was about 0.03 ms. The results show furthermore that the XORing at the client took about 0.18 ms on the average. It appears slightly slower for GF (2^{16}). This reinforces the similar findings in Section 5.2.

The individual insert time for $k = 0$ is about 15% longer than that of an individual update. This is the price for the internal LH splits within the bucket for the inserts. Next, if the updates used 2PC, then the data bucket would need to issue an additional k commit messages. An update time for, for example, $k = 3$ and GF (2^{16}), assuming finally 0.03 ms per message would need at the data bucket about $0.59 + 0.09 = 0.68$ ms. Hence the throughput would be almost 1500 updates/s. This is less than for the inserts, but the number still appears rather very attractive by comparison to the present disk files even for $k \leq 1$ only.

5.7 Bucket Recovery

The recovery manager performs this manipulation as in Section 4.1. For implementation related reasons however, our prototype locates the recovery manager at a parity bucket and not at a spare. To measure the performance, we simulated the creation of an LH_{RS}* group with 4 data buckets and 1, 2, or 3 parity buckets. The group contained $125,000 = 4 \cdot 31,250$ data records consisting again of a 4 B key and 100 B nonkey data. We then reconstructed one, two, and three buckets, made unavailable. We neglected the synchronization phase, whose influence on the response time would be minimal anyhow. The recovery manager loops conceptually over all the existing record groups, over all the parity records in the parity bucket (Section 4.1). In fact, it recovers records by *slices* of a given size s . It requests s successive records from each of the m data/parity buckets, and recovers the s record groups. Then, it requests the next s records from each bucket. While waiting, it sends the recovered slice to the spare(s). Figure 8 presents the effect of slice size on the recovery of a data bucket in the sample case of using the first parity bucket with 1's only and GF(2^{16}). Since the operation is much longer than those of individual records discussed till now, we measured not only the total time (T), but also the process time (P), and the communication time (C).

We determined that the recovery time greatly decreases for a larger s . For $s = 1$, we have $C = 149$ s, $P = 1.735$ s and $T = 165$ s. Figure 8 does not give these values since they are so large, but rather displays values only for $s \geq 100$. Once s is above 1000, T drops under 1s, and P and C under 0.5 s. All the times decrease slightly for larger s and become constant when we choose s over 3000. This is a consequence of our latest communication architecture based on the already mentioned passive TCP connections. The result means that a server may efficiently work with buffers much smaller than the bucket capacity b , for example, 10 times smaller. The experiments with our earlier architectures are in Moussa [2003]. They show the clear superiority of our current implementation.

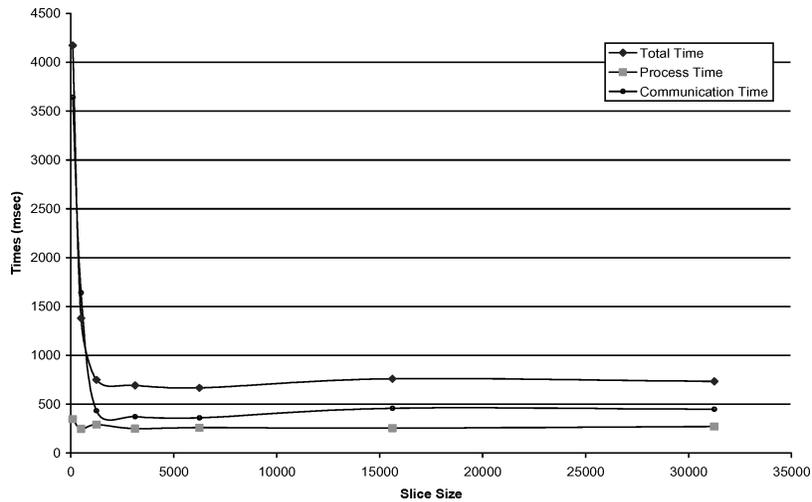
Fig. 8. A single data bucket recovery time (milliseconds) as function of the slice size s .

Table II. Best Data Bucket Recovery Times (seconds) and Slice Sizes

	$GF(2^8)$				$GF(2^{16})$			
	s	T	P	C	s	T	P	C
1-DB (XOR)	15625	0,520	0,225	0,296	6250	0,552	0,240	0,312
1-DB (RS)	6250	0,932	0,630	0,297	15625	0,656	0,354	0,297
2-DBs	15625	1,464	1,156	0,302	15625	0,875	0,562	0,281
3-DBs	6250	2,094	1,711	0,374	15625	1,188	0,823	0,361

Table II completes Figure 8 by listing the T , P , C times for s values minimizing T and $k = 1, 2, 3$. We used $GF(2^8)$ and $GF(2^{16})$. The difference between a T value and the related value of $P + C$ is thread synchronization and switching time. We have also measured all these times for the other s values in Figure 8. For $s \geq 1250$, the differences from the times listed here were under 15% for 1-DB (data bucket) recovery, 5% for 2-DB recovery, and 2% for 3-DBs. The first line of the table presents the 1-DB recovery using the XOR decoding only, as in Figure 8. The second line shows 1-DB recovery using the RS decoding (with the XORing and multiplications). We used another parity bucket instead of the one with ones only, just as in Litwin and Schwarz [2000]. The XOR calculus showed itself to be notably faster for both GF 's used. The gain was expected, but not its actual magnitude. P becomes indeed almost three times smaller for $GF(2^8)$, and almost 1.5 times smaller for $GF(2^{16})$. T decreases less, as it contains the C value. This value is naturally rather stable and turns out to be relatively important with respect to P , despite our fast 1 Gbps network. For the RS decoding we have $C > 0.5P$ at least. Even more interestingly, we reach $C > P$ for the XOR decoding.

All together, our numbers prove the efficiency of the LH_{RS}^* bucket recovery mechanism. It takes only 0.520 s to recover 1 DB in our experiments, and less than 1.2 s to recover 3 DBs, 9.75 MB of data in three buckets. The growth of T appears to be sublinear with respect to the number of buckets recovered. This

Table III. Parity Bucket Recovery Times (seconds) for the Slice
Size of $s = 31\,250$ Records

	$GF(2^8)$			$GF(2^{16})$		
	T	P	C	T	P	C
PB (XOR)	1.872	1.316	0.317	2.062	1.484	0.322
PB (RS)	2.228	1.656	0.307	2.103	1.531	0.322

is a consequence of parallelism at the implementation level, and of the recovery of a bucket group as a whole at the conceptual level. The numbers convincingly confirm the advantage of using $GF(2^{16})$. It halves P of any recovery measured, with the exception of recovery using XOR only. This was the rationale for our choice of this field for the basic LH_{RS}* scheme, given also that with it, parity calculation is as fast as using $GF(2^8)$. Notice that C in Table II increases more moderately than T as the function of the number of DBs recovered.

As discussed in Section 3.4, we used the logarithms of the coefficients in \mathbf{H}^{-1} to obtain the recovery times in Table II. We also experimented using \mathbf{H}^{-1} directly. The results for P were slower, up to 10% for 3-DB recovery, confirming our choice.

The flat character of charts in Figure 8 for larger values of s confirms the scalability of the scheme. It allows us to also guess the recovery times for larger buckets. We can infer from the above numbers that we recover a data bucket group of size $m = 4$ from 1-unavailability at the rate (speed) of 5.89 MB/sec of data. Next, we recover two data buckets of the group at the rate of 7.43 MB/sec. Finally, we recover the group from 3-unavailability at the rate of 8.21 MB/sec. If we have thus, for instance, 1 GB of data per bucket, the figures imply a value for T of about 170 sec for 1-DB recovery, 270 sec for 2 GB recovered, and 365 sec, about 6 min per 3 GB recovered, respectively. If we choose the group size $m = 8$, to halve the storage overhead, the recovery rates will halve as well, while the recovery time will double, and so on.

Table III presents the parity bucket recovery time, again for 31,250 records to recover and $s = 31,250$. The time T to recover bucket P0 using XOR only, analyzed in the row marked PB (XOR), is faster than for the other buckets using the RS calculus. We again observe fast performance. The XOR-only recovery using 2 B symbols (from $GF(2^{16})$) is less efficient than that using 1B symbols (from $GF(2^8)$). We have a reverse picture for the other parity bucket, as the last row in Table III shows. The small difference in P value with respect to the one reported in Section 5.2 is due to the experimental nature of the analysis. The measurements naturally vary slightly among experiments. Similarly as for data buckets, Table III allows us to infer parity bucket recovery rates per MB of data stored for various values of m and the recovery times of the parity buckets of various sizes.

5.8 Record Recovery

Our prototype places the record recovery manager at one of the parity buckets. It acts as described in Section 4.2. Table IV shows the average total record recovery time T we measured. The bucket size was $b = 50,000$. The group size

Table IV. Record Recovery Times (milliseconds)

$GF(2^8)$		$GF(2^{16})$	
XOR	RS	XOR	RS
1.285	1.308	1.297	1.327

was again $m = 4$. The times are measured at the parity bucket. They start when the bucket gets the message from the coordinator, and end with the recovery of the record.

The times for $GF(2^{16})$ are slightly higher. The reason is that we convert 1B characters to 2B symbols and back. In any case, we measured the average scan time of our parity bucket to locate the key c of the data record, as described in Section 4.2, to be 0.822 ms. This is the dominant part of the total time (62% and 64% respectively).

The results match our intuition and the experimental key search times. They confirm that the basic record recovery capability should prove to be often sufficient in practice. The deterioration of the search time with respect to the normal one, we recall of 0.24 ms, is nevertheless about 5.5 times. If one seeks faster record recovery, or if the buckets are much larger, the additional index (c, r) at the parity bucket, mentioned in Section 4.2, should help. For $GF(2^{16})$ and our first parity bucket, one may estimate the decrease to almost $1.296 - 0.822 = 0.474$ ms. The ratio to the normal time becomes less than twice. Notice however the price tag for the index: more storage at the parity bucket and additional processing of an insert and delete at the parity bucket. Knowledge of the scan time allows us to further evaluate the record recovery time for other values of m or b . The communication and processing times are about linear with m , while the bucket scan time is linear with b . Notice finally that even the basic record recovery times remain significantly faster than for a disk file. In our case, the typical ratio should be at least about eight times.

6. RELATED WORK

Traditionally, high availability was not part of a (key-based) data structure in both centralized and distributed environments. If needed, a lower storage level such as mirroring or RAID-like techniques provided it. This approach simplifies the design of a data structure, but it can deteriorate access times in a distributed environment. For example, a dictionary data structure using hashing could place a data unit at some particular node. However, the underlying RAID system could move the data to a different node or even distribute it over several nodes. This lower-level interference would result in additional messaging that an integration of the parity data management into the hashing structure could avoid.

The problem is more acute for a scalable distributed storage environment with a large number of nodes. The elementary reliability calculus shows that higher levels of availability are often necessary for a data structure stored on many nodes. One approach provides the high level at each node. This approach fails if the storage nodes are standard PCs or workstations, especially in a P2P network where nodes may have low availability [Weatherspoon and Kubiatowicz 2002]. In addition, files in the same environment may require

different availability levels just because of their different sizes. The alternative is to integrate high availability into scalable distributed data structures and let the availability level itself scale.

The concept of a *high-availability data structure* appeared in response to the need for integrating high-availability and SDDS [Litwin and Neimat 1996]. The first high-availability SDDS was LH_M*, where high-availability results from mirroring two LH* files. The files contain exactly the same records. They may however differ in their internal structures, for example, the bucket size. In any case, the two files in LH_M* are more strongly coupled than usual mirrors. LH_M* can even recover some cases of double or more unavailability.

Litwin et al. [1997] proposed another 1-availability SDDS called LH_S*. LH_S* partitions a record into n segments, stored at n different sites. It adds an $(n+1)^{st}$ XOR parity segment at some other site. Compared with LH_M*, the parity overhead is much smaller, close to $1/n$. Operations require, in contrast, more messages. An, LH_S* key search in normal mode needs n messages, even though the messages are shorter.

Another 1-available SDDS, LH_g*, [Litwin and Risch 1997, 2001; Lindberg 1997] keeps records intact. It introduces the concept of record groups used by LH_{RS}*. Retrospectively, the LH_{RS}* parity calculus generalizes LH_g* to higher availability. As for LH_{RS}*, an LH_g* record enters a record group when it is created. The group members are always on different servers and the group contains an additional parity record of the same structure as a LH_{RS}* parity record. The initial record group is the same for an LH_g* record as for an LH_{RS}* record. However, an LH_g* record keeps its initial record group membership, regardless of its moves caused by splits. In comparison to LH_{RS}*, LH_g* splits are faster. In contrast, a data bucket recovery processing is more costly. In particular, one always scans all the parity buckets, instead of usually only one for LH_{RS}*. Notice that the recovery is not then necessarily longer than for LH_{RS}*, as the scans can be parallel. If the communication is slow with respect to the processing time, it can be even faster.

LH_{SA}* was the first SDDS to achieve scalable availability [Litwin et al. 1998, 1999]. To achieve k -availability, LH_{SA}* places each record in k or $k+1$ different record groups that only intersect in this one record. Each record group has an additional parity record, basically consisting of the XOR of the other records in the group. LH_{SA}* places the buckets conceptually into a high-dimensional cube with n buckets in the first k or $k+1$ dimensions. Just as for LH_{RS}*, a controlled or an uncontrolled strategy adds parity buckets. A small LH_{SA}* file with a $k > 1$ has a larger storage overhead than a corresponding LH_{RS}* file. This advantage of LH_{RS}* dissipates for larger files. LH_{SA}* parity calculations use only XORing, which gives it an advantage over k -available LH_{RS}* files for $k > 1$. However, if there is more than one unavailable bucket, recovering a lost record can involve additional recovery steps. A deeper comparison of trade-offs between LH_{SA}* and LH_{RS}* remains to be undertaken.

Outside the domain of SDDSs, research has addressed high-availability needs for distributed flat files for many years. The dominant approach was replication [Haskin and Schmuck 1996]. The major issue was consistency of replicas [Paris 1993]. Disk arrays in a centralized environment historically needed high

availability with less storage overhead [Burkhard and Menon 1993; Hellerstein et al. 1994]. The arrays typically have a fixed number of disks so that the proposed high-availability schemes were static. The aspects under investigation were mainly the parity update mechanisms (e.g. parity logging), and the parity placement providing the 1-availability through XORing. These were the performance determinants of a disk array. Next, parity placement schemes appeared to be intended for larger, but still static, arrays, for example, Alvarez et al. [1997]. Current research increasingly focuses on very large storage systems, using an expandable number of storage units, whether disks or entire servers. Recent proposals for the k -available ($k > 1$) erasure correcting codes discussed in App. C (Online Appendix available in the ACM Digital Library) came from this context.

High-availability is also a general goal for a DBMS. Nevertheless, our aspect of this concept, the unavailability of a part of data storage, has received relatively little attention. The general assumption seems to be the use of a high-availability storage or file system underneath. Typically, it should be software or hardware RAID storage. For a parallel DBMS, this should concern each DBMS node. At the database layer, replication seems the only technique used. The DBMS is then typically 1-available, with respect to storage node unavailability.

The Clustra DBMS, now a commercial product, proposes a DBMS level structure that some claim to be the most efficient in the domain [Sabaratnam et al. 1999]. It hashes partitions of a table into fragments located each at a different node. The nodes communicate using a dedicated high-speed switch. Clustra hashing is static, hence has limited scalability compared to LH_{RS}^* . The practical limit is 24 nodes at present. Each fragment is replicated on two nodes, using the primary copy approach. If a fragment is unavailable, (detected by lack of heart beat), its available copy, possibly the primary one, is copied to a spare. The partitioning typically limits the recovery to a single fragment. The whole scheme makes Clustra tables only one-available and limits their scalability compared to our scheme. This conclusion holds for other prominent DBMSs, whether they use for the parallel table partitioning, the (static) hashing (DB2), or range partitioning (SQL Server) or both (Oracle).

Research has also started addressing the high-availability needs of scalable disk farms [Xin et al. 2003, 2004]. These should be soon necessary for grid computing and very large Internet databases. Some simple techniques are already in everyday use. They are apparently replication based, but covered by corporate secrecy. The prominent example is Google. The gray literature estimates its farm spreading already over more than 10,000 Linux nodes, perhaps as many as 54,000 [Donoghue 2003; Economist 2003]. There are also open research proposals for high-availability distributed data structures over large clusters specifically intended for the Internet access. One is a distributed hash table scheme with built-in specific replication [Gribble et al. 2000]. An ongoing research project follows up with the goal of a scalable distributed highly-available linked B-tree [Boxwood 2003].

Emerging P2P applications, including those based on Wi-Fi, also have compelling high-availability storage needs [Anderson and Kubiatowicz 2002;

Kubiatowicz 2003; Dingledine et al. 2000]. In this new environment the availability of the nodes should be more “chaotic” than one typically supposed in the past. Their number and geographical dispersion can also be larger by orders of magnitude, possibly running into hundreds of thousands in the near future and later reaching millions, spread worldwide. This thinking clearly shares some rationales for LH_{RS}*. Our scheme could thus turn out to be useful for these new applications as well.

7. CONCLUSION

LH_{RS}* is a high-availability scalable distributed data structure. It scales up to any size and any availability level k that one can reasonably foresee for an application these days. File scalability is transparent to the application, as for any SDDS. The k -availability may scale transparently as well, or may be adjusted by the application on demand.

The scheme matured in many aspects from our initial proposal [Litwin and Schwarz 2000]. Changes were made to the parity calculus and to various algorithmic issues to make the file always at least $(K - 1)$ -available and make the parity calculus as fast as possible. We have thus increased the Galois Field size to $GF(2^{16})$. We have changed the parity matrix \mathbf{P} so it has a first column and first row of ones. We have also improved the calculus by taking advantage of logarithmic parity. We have built a prototype implementation that proves the feasibility of the scheme. We have experimented on this basis with the new, and the former, parity calculus, as well as with the previously mentioned algorithmic issues. Performance analysis showed a substantial speedup of various operations.

At present, for the most frequent case of $k = 1$, our erasure correction uses XORing only and is thus as fast as possible to our best knowledge, under our constraints, especially that of storage efficiency. For $k > 1$, it appears more effective in practice than if we used any alternative parity code or scheme we are aware of. This is also true for our own earlier approach, as we just mentioned. The as yet unique presence of the row of ones contributes to this performance. Likewise, our original use of the logarithmic matrices. In particular, while the parity storage and communication overheads increase substantially with the k used, they globally always remain close to the optimal bounds. Another known high-availability SDDS scheme may nevertheless eventually outperform the LH_{RS}* on a selected feature. The diversity should profit any application.

Finally, our prototype implementation has shown very fast access and recovery performance. Our test-bed files with 125K records recovered in less than a second from a single unavailability and in about two seconds from a triple one. Individual search, insert, and update times were at most 0.5 msec for a 3-available file. Bulk operations were many times faster. This performance is partially due to processing data in the distributed RAM. Altogether, the capabilities of our scheme should attract numerous applications of hash files, for which LH_{RS}* should be “plug-compatible,” as we discussed in the introduction. To recall, these applications are potentially very numerous, as hash files are ubiquitous. This area includes new domains of grid computing and of P2P, and

popular DBMSs. The latter still use the more limited static and 1-available replication or RAID storage for high-availability.

Future work should concern experiments with applications of our scheme. One should also port the parity subsystem to other known 0-available SDDS schemes. The range partitioning schemes appear to be preferred candidates. One should also add the capabilities of concurrent and transactional access to LH_{RS}^* . Notice that the data records of a record group conflict on the parity records. One should finally study in more depth the discussed variants, including those in Online Appendix C available in the ACM Digital Library.

APPENDIX

A. PARITY MATRICES

We present the first 32 row by 10 column submatrices of the generic parity matrix \mathbf{P}' and of the generic logarithmic parity matrix \mathbf{Q}' for $GF(2^{16})$ that we use for LH_{RS}^* . The values are four hexadecimal digits. The submatrices allow for actual matrices \mathbf{P} and \mathbf{Q} for groups of size m up to 32, with k up to 10, values that should suffice in practice. Next, we show 32×20 portions of \mathbf{P}' and \mathbf{Q}' for $GF(2^8)$ used in the examples. The entries of \mathbf{P}' are now $GF(2^8)$ elements given as two hexadecimal digits. The entries of \mathbf{Q}' are logarithms, given as decimal numbers between 0 and 254. The program to generate the complete matrices can be requested from the authors at CERIA [<http://ceria.douphine.fr>].

0001	0001	0001	0001	0001	0001	0001	0001	0001	0001
0001	eb9b	2284	9e44	f91c	7ab9	2897	41f6	a9dd	5933
0001	2284	9e74	d7f1	0fe3	79bb	5658	efa6	30f3	641c
0001	9e44	d7f1	75ee	512d	4e14	16bb	2ce0	36c8	0f9a
0001	f91c	0fe3	512d	59c3	d037	b205	cb3c	f6e2	c606
0001	7ab9	79bb	4e14	d037	b259	e9b9	2c40	81b2	70b5
0001	2897	5658	16bb	b205	e9b9	c7b6	07c7	8670	86ac
0001	41f6	efa6	2ce0	cb3c	2c40	07c7	2c2c	5ddc	148f
0001	a9dd	30f3	36c8	f6e2	81b2	8670	5ddc	7702	1f19
0001	5933	641c	0f9a	c606	70b5	86ac	148f	1f19	9c98
0001	52d5	59c3	94f7	4d4d	e9b2	40f1	2d00	9c04	bc3f
0001	3f68	3d2b	00f1	32c2	dfb2	6ab4	c6a6	9eba	0241
0001	47b5	cb3c	6f1f	2d00	39e6	799c	c83b	92df	c24d
0001	8656	b46f	59c3	903a	7432	9aef	46f5	3b50	e867
0001	db71	b612	eb07	496e	ac26	74c2	04cc	5f5d	23b9
0001	92fe	acb3	3045	fef2	7607	ad10	2df0	0b2f	1eaa
0001	99f1	6d93	5803	1ce8	4099	136c	af32	35b6	3274
0001	2c68	2d00	4fe4	50bf	16af	88ec	2ec0	bfa2	2b90
0001	7502	c6a6	8f58	5a03	2887	89ca	6724	e0be	39e1
0001	227d	16a8	eacf	0f59	67af	7702	838d	3517	85a1
0001	1027	496e	a06a	c486	61ab	131f	2e00	b405	fafc
0001	0a46	87b4	74c2	f85b	e8ef	5b03	3163	8b11	b4a5
0001	27df	9523	379a	7e8d	0301	0221	7702	fe93	d06b
0001	3dad	27bc	9f64	7602	5cf1	ee22	2e48	64a0	a751
0001	8dd7	f6e2	f125	9c04	f720	c0a3	92df	ee07	1d73
0001	1e27	fd33	93fd	86cd	9ca4	2614	9961	8483	c26e
0001	0cf7	46f5	2d00	cf2f	6f7b	2f80	8497	4baf	2900
0001	651e	a0d6	58d3	dbb0	96ed	f57a	2f20	5c03	69d3
0001	71ec	0f8d	496e	7702	51a4	f85b	ba6f	a34b	ce4b
0001	00a1	fd45	564f	dc72	6924	b699	9d44	de6d	a90c
0001	b4ca	b385	025b	0eb8	47bd	31f6	f173	a768	798b
0001	59c3	73b6	b325	4b4b	6ba7	7ca5	2fd0	5e55	9ac4

Fig. 9. Generic parity matrix \mathbf{P}' for $GF(2^{16})$: first 32 rows by 10 columns.

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	105	200	55	12	243	86	169	29	226	240	15	107	148	230	25	62	193	34	221
0	120	185	135	41	115	214	140	136	92	119	163	50	173	205	82	74	46	181	209
0	175	80	240	127	29	226	128	121	107	148	134	25	198	57	230	108	12	243	147
0	239	93	234	133	16	162	21	54	132	223	213	201	123	32	42	189	236	228	250
0	227	179	148	28	121	107	76	161	25	198	238	230	94	17	57	43	127	29	194
0	99	52	119	203	136	92	156	104	50	173	40	82	215	151	205	47	41	115	190
0	13	64	107	148	191	242	6	79	75	25	188	67	230	180	176	103	240	224	81
0	38	225	10	91	186	100	244	96	217	30	245	164	69	155	11	210	97	7	157
0	9	210	25	198	79	75	14	246	67	230	45	176	57	241	180	159	148	191	228
0	130	89	173	141	104	50	71	166	82	215	125	205	184	114	151	81	203	136	51
0	145	118	144	116	129	157	219	111	137	110	230	36	98	126	139	10	19	85	102
0	170	193	223	37	54	132	31	218	201	123	224	42	85	62	32	192	133	16	73
0	63	168	248	66	25	117	46	230	189	209	138	192	190	7	87	195	130	107	237
0	88	143	50	173	146	251	194	4	61	82	109	112	205	247	167	90	119	118	87
0	113	250	198	158	161	25	165	90	230	94	97	57	5	142	17	181	28	121	84
0	202	39	128	102	199	174	144	86	6	76	20	14	165	117	233	220	53	216	127
0	140	73	94	211	90	230	88	35	57	5	91	17	162	208	142	115	158	161	182
0	248	220	82	215	4	61	84	212	112	205	169	167	151	219	247	35	173	146	7
0	214	27	20	16	203	170	230	141	183	154	220	3	60	222	244	235	228	41	112
0	183	12	123	168	218	201	149	68	42	85	191	32	129	108	62	87	37	54	106
0	74	68	67	230	156	235	115	71	39	176	100	236	180	37	133	99	25	140	20
0	243	154	63	242	9	127	103	221	28	187	205	158	235	237	211	128	152	13	246
0	187	8	209	21	230	189	41	57	192	190	202	87	51	186	7	214	66	25	234
0	44	46	30	105	96	217	222	131	164	69	118	11	13	74	155	124	91	186	137
0	95	230	101	108	93	53	131	193	102	103	84	120	159	130	99	153	62	171	152
0	150	175	179	3	82	64	236	205	210	250	72	124	73	216	95	5	183	50	45
0	221	124	230	94	246	67	233	239	176	57	10	180	17	70	241	245	198	79	16
0	80	55	201	123	60	208	51	65	145	42	113	77	32	101	160	178	223	154	95
0	77	146	110	72	111	137	122	174	36	98	57	139	225	135	126	30	116	129	120
0	66	157	215	249	166	82	200	178	205	184	53	151	78	27	114	228	141	104	177
0	230	160	212	178	237	31	251	234	149	38	199	185	44	89	52	203	166	211	70

Fig. 12. Generic logarithmic parity matrix \mathbf{Q} for $GF(2^8)$: first 32 rows by 20 columns.

B. DEFINITION OF TERMS

Term	Description	Typical Value
	Addressing	
F	an LH_{RS}^i file	initially 0, scales monotonically
i	file level	
n	split pointer	$0-2'-1$
(i, n)	file state	$0-i$
N	current number of data buckets in the file	$N = 2' + n$
a	logical address of a data bucket server	$[0; 1; 2; \dots; M-1]$
A	physical address of a data bucket server	IP address
A_0	initial physical address of the file (server of data bucket 0)	IP address
j	data bucket level	i or $i + 1$
c	(primary) key of a data record	random; $0-2^{32} - 1$
h_i	series of hash functions	$C \bmod N * 2'$
α	load factor	$0.6-1.0$
	Parity calculus	
$GF(2^f)$	Galois Field of size (2^f)	$GF(2^{16})$
F	Galois Field of size (2^8)	
ECC	Erasur Correcting Code	
RS	Reed-Solomon Code	
P	parity field (in parity record)	$GF(2^{16})$ symbols
C	record group structure field (in parity record)	c_0, c_1, \dots, c_{m-1}
k	bucket (record) group local availability level	$1-10$

K_{file}	global file availability level	1—10
g	bucket group number	1, 2 . . .
r	data record rank (and parity record key)	1— b
α	primitive element in GF	$\alpha = 2$
$\log_{\alpha}(\zeta)$	logarithm of symbol ζ ; $\zeta \in GF(2^f)$, $\zeta \neq 0$	Table 5
antilog(f)	antilogarithm of integer i ; $0 \leq i < 2^f - 1$	Table 5
I	identity matrix $m \times m$	
P'	generic parity matrix	Figure 9
P	actual parity matrix	upper left $m \times K$ submatrix of P'
Q'	generic logarithmic parity matrix	Figure
Q	actual logarithmic parity matrix	upper left $m \times K$ submatrix of Q'
H, H⁻¹	decoding matrices ($m \times m$, formed from avail. columns of P)	
L_A	list of available buckets in a bucket group recovery	m
L_S	list of spare buckets for a bucket group recovery	$l \leq k$
File Param.	Description	Typical Value
B	bucket capacity (records per data bucket)	50—1,000,000
K	intended file availability level (scales monotonically)	1—5
m	bucket group size (also max. record group size)	4—32

ACKNOWLEDGMENTS

Many thanks to Jim Gray and to the anonymous referees for helpful suggestions.

REFERENCES

- ALVAREZ, G., BURKHARD, W., AND CRISTIAN, F. 1997. Tolerating multiple failures in RAID Architecture with Optimal Storage and Uniform Declustering. In *International Symposium on Computer Architecture, ISCA-97*, 62–72.
- ANDERSON, D. AND KUBIATOWICZ, J. 2002. The Worldwide Computer. In *Scientific American* 286, 3, March.
- THE BOXWOOD PROJECT. <http://research.microsoft.com/research/sv/Boxwood/>.
- BARTALOS, G. 1999. Internet: D-day at eBay. *Yahoo INDIVIDUAL INVESTOR ONLINE*, (July 19).
- BERTINO, E., OOI, B. C., SACKS-DAVIS, R., TAN, K. L., ZOBEL, J., SHIDLOVSKY, B., AND CATANIA, B. 1999. Indexing Techniques for Advanced Database Systems. Kluwer.
- BENNOUR, F., DIÈNE, A., NDIAYE, Y., AND LITWIN, W. 2000. Scalable and distributed linear hashing LH_{LH}* under Windows NT. In *SCI-2000 (Systemics, Cybernetics, and Informatics)*, Orlando, Florida.
- BENNOUR, F. 2002. Performance of the SDDS LH_{LH}* under SDDS-2000. In *Distributed Data and Structures 4 (Proceedings of WDAS 2002)*, Carleton Scientific, 1–12.
- BURKHARD, W. A. AND MENON, J. 1993. Disk array storage system reliability. In *Proceedings of the 22nd International Symposium on Fault Tolerant Computing*, Toulouse, 432–441.
- BEN-GAN, I., AND MOREAU, T. *Advanced Transact-SQL For SQL Server 2000*. 2003. Apress, ISBN 1-8931115-82-8.
- BREITBART, Y., VINGRALEK, R., AND WEIKUM, G. 1996. Load control in scalable distributed file structures. *Distributed and Parallel Databases* 4, 4, 319–354.
- BREITBART, Y. AND VINGRALEK, R. 1998. Addressing and balancing issues in distributed B+ trees. In *1st Workshop on Distributed Data and Structures (WDAS '98)*, Carleton-Scientific.
- COM. ACM. 1997. Special Issue on High-Performance Computing (Oct).
- CERIA Home page: <http://ceria.dauphine.fr/>
- <http://www.contingencyplanningresearch.com/cod.htm>. 1996. Cost of a downtime Study.
- DINGLEDDINE, R., FREEDMAN, M., AND MOLNAR, D. 2000. The free haven project: Distributed anonymous storage service. *Workshop on Design Issues in Anonymity and Unobservability* (July).

- DONOGHUE, A. BOLDLY 2003. Googling into the future. <http://insight.zdnet.co.uk/internet/ecommerce/0,39020454,39116781,00.htm>
- ECONOMIST 2003. Moving up the stack. www.economist.com. May.
- GRIBBLE, S., BREWER, E., A., HELLERSTEIN, J., AND CULLER, D. 2000. Scalable, Distributed Data Structures for Internet Service Construction. 4th Symposium on Operating Systems Design and Implementation (OSDI'00).
- GRAY, J., SZALAY, A. S., IHAKAR, A., KUNSZT, P. S., STOUGHTON, C., SLUTZ, D. R., AND VAN DEN BERG, J. 2002. Data Mining of SDDS SkyServer Database. International Workshop on Distributed Data Structures, (WDAS'02), Carleton Scientific.
- HASKIN, R. AND SCHMUCK, F. 1996. The Tiger Shark File System. COMPCON-96, 1996.
- HELLERSTEIN, L., GIBSON, G., KARP, R., KATZ, R., AND PATTERSON, R. 1994. Coding techniques for handling failures in large disk arrays. *Algorithmica*, vol. 12, p. 182–208.
- KNUTH, D. 1998. *The art of computer programming. Vol. 3 Sorting and searching*. 2nd Ed. Addison-Wesley, 780.
- KUBIATOWICZ, J. 2003. Extracting guarantees from chaos. In *Communications of the ACM*, 46, 2, Feb.
- KARLSON, J., LITWIN, W., AND RISCH, T. 1996. LH_{LH}^* : A scalable high performance data structure for switched multicomputers. In Apers, P., Gardarin, G., Bouzeghoub, M., (eds.) *Extending Database Technology*, EDBT96, Lecture Notes in Computer Science, vol. 1057. Springer Verlag.
- LINDBERG, R. 1997. A Java Implementation of a Highly Available Scalable and Distributed Data Structure LH^*g . Master Th. LiTH-IDA-Ex-97/65. U. Linkoping, 1997/62.
- LITWIN, W. 1994. Linear hashing: A new tool for file and table addressing. Reprinted from *VLDB80* in *Readings in Databases*, edited by M. Stonebraker, 2nd Edition, Morgan Kaufmann Publishers.
- LITWIN, W. 1980. Linear hashing: A new algorithm for files and tables addressing. International Conference on Databases. Aberdeen, Heyden, p. 260–275.
- LITWIN, W., NEIMAT, M.-A., LEVY, G., NDIAYE, S., AND SECK, T. 1997. LH_S^* : A high-availability and high-security Scalable Distributed Data Structure. IEEE-Res. Issues in Data Eng. (RIDE-97).
- LITWIN, W., MENON J., AND RISCH, T. 1998. LH^* with Scalable Availability. IBM Almaden Res. Rep. RJ 10121 (91937), (May).
- LITWIN, W., MENON, J., RISCH, T., AND SCHWARZ, T. 1999. Design Issues For Scalable Availability LH^* Schemes with Record Grouping. DIMACS Workshop on Distributed Data and Structures, Princeton U. Carleton Scientific.
- LITWIN, W., MOUSSA, R., AND SCHWARZ, T., 2004a. LH_{RS}^* : A Highly Available Distributed Data Storage System. Research Prototype Demonstration. VLDB Toronto.
- LITWIN, W., MOUSSA, R., AND SCHWARZ, T. 2004b. LH_{RS}^* : A Highly Available Distributed Data Storage System. CERIA Tech. Rep. (Dec).
- LITWIN, W., NEIMAT, M.-A., AND SCHNEIDER, D. 1993. Linear Hashing for Distributed Files. ACM-SIGMOD International Conference on Management of Data.
- LITWIN, W., NEIMAT, M.-A., AND SCHNEIDER, D. 1996. A Scalable Distributed Data Structure. *ACM Trans. Datab. Syst.*, Dec.
- LITWIN, W. AND NEIMAT, M.-A. 1996. High-Availability LH^* Schemes with Mirroring, International Conference on Cooperating Information Systems, (COOPIS) IEEE Press.
- LITWIN, W. AND RISCH, T. 1997. LH^*g : A High-availability Scalable Distributed Data Structure through Record Grouping. Res. Rep. CERIA, U. Dauphine and U. Linkoping (May).
- LITWIN, W. AND RISCH, T. 2001. LH^*g : A high-availability scalable distributed data structure by record grouping. *IEEE Trans. Knowl. Data Eng.* 14, 4, 923–927.
- LITWIN, W. AND SCHWARZ T. 2000. LH_{RS}^* : A high-availability scalable distributed data structure using Reed Solomon codes. ACM-SIGMOD International conference on Management of Data.
- LITWIN, W. AND RISCH, T. 2002. LH^*g : A High-availability scalable distributed data structure by record grouping. *IEEE Trans. Knowl. Data Eng.* 14, 4, July/Aug.
- LJUNGSTRÖM, M. 2000. Implementing LH_{RS}^* : A scalable distributed highly-available data structure, Master Thesis, Feb., CS Dep. U. Linkoping, Sweden.
- LUBY, M., MITZENMACHER, M., SHOKROLLAHI, M., SPIELMAN, D., AND STEMANN, V. 1997. Practical Loss-Resilient Codes, STOC 97, Proceedings of the twenty-ninth annual ACM Symposium on Theory of Computing, El Paso, TX, 150–159.

- MACWILLIAMS, F. J. AND SLOANE, N. J. A. 1997. *The Theory of Error Correcting Codes*. Elsevier/North Holland, Amsterdam.
- MOUSSA, R. 2003. In *Distributed Data and Structures 4*, Carleton Scientific (Records of WDAS 2002, Paris).
- MOUSSA, R. 2004. Experimental Performance Analysis of LH_{RS}*. CERIA Res. Rep. [CERIA].
- MOUSSA, R. AND LITWIN, W. 2002. Experimental performance analysis of LH_{RS}* parity management. *Distributed Data and Structures 4, Records of the 4th International Meeting (WDAS 2002)*, Paris, France.
- PÁRIS, J. F. 1993. The management of replicated data. In *Proceedings of the Workshop on Hardware and Software Architectures for Fault Tolerance*. Mt. St. Michel, Fr. June.
- RAMAKRISHNAN, R. 1999. *Database Management Systems*. McGraw Hill.
- RFC 793—Transmission Control Protocol <http://www.faqs.org/rfcs/rfc793.html>
- SABARATNAM M., TORBJORNSEN, AND HVASSHOVD, S.-O. 1999. Evaluating the effectiveness of fault tolerance in replicated database management systems. *29th. Annual International Symposium on Fault Tolerant Computing*.
- SCHWARZ, T. 2003. Generalized Reed Solomon Codes for Erasure Correction in SDDS. *Workshop on Distributed Data and Structure 4, WDAS-4, Paris. Carleton Scientific*.
- SDDS-BIBLIOGRAPHY. <http://192.134.119.81/SDDS-bibliographie.html>, <http://ceria.dauphine.fr/witold.html>
- VINGRALEK, R., BREITBART, Y., WEIKUM, G. AND SNOWBALL. 1998. Scalable storage on networks of workstations. *Distributed and Parallel Databases 6, 2*, 117–156.
- WEATHERSPOON, H. AND KUBIATOWICZ, J. 2002. Erasure coding vs. replication: A quantitative comparison. *1st International Workshop on Peer-to-Peer systems, IPTPS-2002*. March.
- XIN, Q., MILLER, E., SCHWARZ, T., BRANDT, S., LONG, D., LITWIN, W. 2003. Reliability mechanisms for very large storage systems. *20th IEEE mass storage systems and technologies (MSST 2003)*, San Diego, CA. 146–156.
- XIN, Q., MILLER, E., AND SCHWARZ, T. 2004. Evaluation of distributed recovery in large-scale storage systems. In *13th IEEE International Symposium on High Performance Distributed Computing (HPDC'04)*, Honolulu, HI.

Received July 2004; revised February 2005; accepted May 2005