

Liquid Cloud Storage

Michael G. Luby, Qualcomm Technologies, Inc. (luby@qti.qualcomm.com)

Roberto Padovani, Qualcomm Technologies, Inc. (padovani@qti.qualcomm.com)

Thomas J. Richardson, Qualcomm Technologies, Inc. (tomr@qti.qualcomm.com)

Lorenz Minder, Qualcomm Technologies, Inc. (lminder@qti.qualcomm.com)

Pooja Aggarwal, Qualcomm Technologies, Inc. (poojaa@qti.qualcomm.com)

A *liquid system* provides durable object storage based on spreading redundantly generated data across a network of hundreds to thousands of potentially unreliable storage nodes. A liquid system uses a combination of a *large code*, *lazy repair*, and a *flow storage organization*. We show that a liquid system can be operated to enable flexible and essentially optimal combinations of storage durability, storage overhead, repair bandwidth usage, and access performance.

CCS Concepts: •**Information systems** → **Distributed storage**; Information storage technologies; Storage recovery strategies; •**Computer systems organization** → **Maintainability and maintenance**; *Reliability*; *Availability*; *Redundancy*;

Additional Key Words and Phrases: distributed information systems, data storage systems, data warehouses, information science, information theory, information entropy, error compensation, time-varying channels, error correction codes, Reed-Solomon codes, network coding, signal to noise ratio, throughput, distributed algorithms, algorithm design and analysis, reliability, reliability engineering, reliability theory, fault tolerance, redundancy, robustness, failure analysis, equipment failure.

ACM Reference Format:

ACM V, N, Article A (January YYYY), 44 pages.

DOI : <http://dx.doi.org/10.1145/0000000.0000000>

1. INTRODUCTION

Distributed storage systems generically consists of a large number of interconnected storage nodes, with each node capable of storing a large quantity of data. The key goals of a distributed storage systems are to store as much source data as possible, to assure a high level of durability of the source data, and to minimize the access time to source data by users or applications.

Many distributed storage systems are built using commodity hardware and managed by complex software, both of which are subject to failure. For example, storage nodes can become unresponsive, sometimes due to software issues where the node can recover after a period of time (transient node failures), and sometimes due to hardware failure in which case the node never recovers (node failures). The paper [8] provides a more complete description of different types of failures.

A *sine qua non* of a distributed storage system is to ensure that source data is durably stored, where durability is often characterized in terms of the *mean time to the loss of any source data (MTTDL)* often quoted in millions of years. To achieve this in the face of component failures a fraction of the raw storage capacity is used to store redundant data. Source data is partitioned into objects, and, for each object, redundant data is generated and stored on the nodes. An often-used trivial redundancy is *triplication* in which three copies of each object are stored on three different nodes. Triplication has several advantages for access and ease of repair but its overhead is quite high: two-thirds of the raw storage capacity is used to store redundant data. Recently, Reed-Solomon (RS) codes have been introduced in production systems to both reduce storage overhead and improve MTTDL compared to triplication, see e.g., [10], [8], and [19]. As explained in [19] the reduction in storage overhead and improved MTTDL achieved by RS codes comes at the cost of a *repair bottleneck* [19] which arises from the significantly increased network traffic needed to repair data lost due to failed nodes.

The repair bottleneck has inspired the design of new types of erasure codes, e.g., *local reconstruction codes* (LR codes) [9], [19], and *regenerating codes* (RG codes) [6], [7], that can reduce repair traffic compared to RS codes. The RS codes, LR codes and RG codes used in practice typically use small values of (n, k, r) , which we call *small codes*. Systems based on small codes, which

we call *small code systems*, therefore spread data for each object across a relatively small number of nodes. This situation necessitates using a *reactive repair* strategy where, in order to achieve a large MTTDL, data lost due to node failures must be quickly recovered. The quick recovery can demand a large amount network bandwidth and the repair bottleneck therefore remains.

These issues and proposed solutions have motivated the search for an understanding of tradeoffs. Tradeoffs between storage overhead and repair traffic to protect individual objects of small code systems are provided in [6], [7]. Tradeoffs between storage overhead and repair traffic that apply to all systems are provided in [14].

Initially motivated to solve the repair bottleneck, we have designed a new class of distributed storage systems that we call *liquid systems*: Liquid systems use *large codes* to spread the data stored for each object across a large number of nodes, and use a *lazy repair strategy* to slowly repair data lost from failed nodes.¹ In this paper we present the design and its most important properties. Some highlights of the results and content of paper include:

- Liquid systems solve the repair bottleneck, with less storage overhead and larger MTTDL than small code systems.
- Liquid systems can be flexibly deployed at different storage overhead and repair traffic operating points.
- Liquid systems avoid transient node failure repair, which is unavoidable for small code systems.
- Liquid systems obviate the need for sector failure data scrubbing, which is necessary for small code systems.
- Liquid system regulated repair provides an extremely large MTTDL even when node failures are essentially adversarial.
- Simulation results demonstrate the benefits of liquid systems compared to small code systems.
- Liquid system prototypes demonstrate improved access speeds compared to small code systems.
- Detailed analyses of practical differences between liquid systems and small code systems are provided.
- Liquid systems performance approach the theoretical limits proved in [14].

2. OVERVIEW

2.1. Objects

We distinguish two quite different types of data objects. *User objects* are organized and managed according to semantics that are meaningful to applications that access and store data in the storage system. *Storage objects* refers to the organization of the source data that is stored, managed and accessed by the storage system. This separation into two types of objects is common in the industry, e.g., see [3]. When user objects (source data) are to be added to the storage system, the user objects are mapped to (often embedded in) storage objects and a record of the mapping between user objects and storage objects is maintained. For example, multiple user objects may be mapped to consecutive portions of a single storage object, in which case each user object constitutes a byte range of the storage object. The management and maintenance of the mapping between user objects and storage objects is outside the scope of this paper.

The primary focus of this work is on storing, managing and accessing storage objects, which hereafter we refer to simply as *objects*. Objects are immutable, i.e., once the data in an object is determined then the data within the object does not change. Appends to objects and deletions of entire objects can be supported. The preferred size of objects can be determined based on system tradeoffs between the number of objects under management and the granularity at which object repair and other system operations are performed. Generally, objects can have about the same size,

¹The term *liquid system* is inspired by how the repair policy operates: instead of reacting and separately repairing chunks of data lost from each failed node, the lost chunks of data from many failed nodes are lumped together to form a “liquid” and repaired as a flow.

or at least a minimum size, e.g., 1 GB. The data size units we use are KB = 1024 bytes, MB = 1024 KB, GB = 1024 MB, TB = 1024 GB, and PB = 1024 TB.

2.2. Storage and repair

The number of nodes in the system at any point in time is denoted by M , and the storage capacity of each node is denoted by S , and thus the overall raw capacity is $\mathcal{D}_{\text{raw}} = M \cdot S$. The storage overhead expressed as a fraction is $\beta = 1 - \frac{\mathcal{D}_{\text{src}}}{\mathcal{D}_{\text{raw}}}$, and expressed as a per cent is $\beta \cdot 100\%$, where \mathcal{D}_{src} is the aggregate size of objects (source data) stored in the system.

When using an (n, k, r) erasure code, each object is segmented into k source fragments and an encoder generates $r = n - k$ repair fragments from the k source fragments. An *Encoding Fragment ID*, or EFI, is used to uniquely identify each fragment of an object, where EFIs $0, 1, \dots, k - 1$ identify the source fragments of an object, and EFIs $k, \dots, n - 1$ identify the repair fragments generated from source fragments. Each of these $n = k + r$ fragments is stored at a different node. An erasure code is MDS (maximum distance separable) if the object can be recovered from any k of the n fragments.

A *placement group* maps a set of n fragments with EFIs $\{0, 1, \dots, n - 1\}$ to a set of n of the M nodes. Each object is assigned to one of the placement groups, and the selected placement group determines where the fragments for that object are to be stored. To ensure that an equal amount of data is stored on each node, each placement group should be assigned an equal amount of object data and each node should accommodate the same number of placement groups.

For the systems we describe, $\frac{\mathcal{D}_{\text{src}}}{\mathcal{D}_{\text{raw}}} = \frac{k}{n}$ and thus $\beta = \frac{r}{n}$, i.e., the amount of source data stored in the system is maximized subject to allotting enough space to store n fragments for each object. Object loss occurs when the number of fragments available for an object is less than k , and thus the object is no longer recoverable. A repairer is an agent that maintains recoverability of objects by reading, regenerating and writing fragments lost due to node failures, to ensure that each object has at least k available fragments. In principle the repairer could regenerate r missing fragments for an object when only k fragments remain available. However, this repair policy results in a small MTTDL, since one additional node failure before regenerating the r missing fragments for the object causes object loss.

A repairer determines the rate at which repair occurs, measured in terms of a *read repair rate*, since data read by a repairer is a fair measure of repair traffic that moves through the network. (The amount of data written by a repairer is essentially the same for all systems, and is generally a small fraction of data read by a repairer.)

Of primary interest is the amount of network bandwidth $\mathcal{R}_{\text{peak}}$ that needs to be dedicated to repair to achieve a given MTTDL for a particular system, as this value of $\mathcal{R}_{\text{peak}}$ determines the network infrastructure needed to support repair for the system. Generally, we set a global upper bound $\mathcal{R}_{\text{peak}}$ on the allowable read repair rate used by a repairer, and determine the MTTDL achieved for this setting of $\mathcal{R}_{\text{peak}}$ by the system. The average read repair rate \mathcal{R}_{avg} is also of some interest, since this determines the average amount of repair traffic moved across the network over time.

The fixed-rate repair policy works in similar ways for small code systems and liquid systems in our simulations: The repairer tries to use a read repair rate up to $\mathcal{R}_{\text{peak}}$ whenever there are any objects to be repaired.

2.3. Small code systems

Replication, where each fragment is a copy of the original object is an example of a trivial MDS erasure code. For example, triplication is a simple $(3, 1, 2)$ MDS erasure code, wherein the object can be recovered from any one of the three copies. Many distributed storage systems use replication.

A Reed-Solomon code (RS code) [2], [17], [12] is an MDS code that is used in a variety of applications and is a popular choice for storage systems. For example, the small code systems described in [8], [10] use a $(9, 6, 3)$ RS code, and those in [19] use a $(14, 10, 4)$ RS code.

When using a (n, k, r) small code for a small code system with M nodes, an important consideration is the number of placement groups P to use. Since $n \ll M$, a large number P of placement groups are typically used to smoothly spread the fragments for the objects across the nodes, e.g., Ceph [22] recommends $P = \frac{100 \cdot M}{n}$. A placement group should avoid mapping fragments to nodes with correlated failures, e.g., to the same rack and more generally to the same failure domain. Pairs of placement groups should avoid mapping fragments to the same pair of nodes. Placement groups are updated as nodes fail and are added. These and other issues make it challenging to design placement groups management for small code systems.

Reactive repair is used for small code systems, i.e., the repairer quickly regenerates fragments as soon as they are lost from a failed node (reading k fragments for each object to regenerate the one missing fragment for that object). This is because, once a fragment is lost for an object due to a node failure, the probability of r additional node failures over a short interval of time when r is small is significant enough that repair needs to be completed as quickly as practical. Thus, the peak read repair rate $\mathcal{R}_{\text{peak}}$ is higher than the average read repair rate \mathcal{R}_{avg} , and \mathcal{R}_{avg} is k times the node failure erasure rate. In general, reactive repair uses large bursts of repair traffic for short periods of time to repair objects as soon as practical after a node storing data for the objects is declared to have permanently failed.

The peak read repair rate $\mathcal{R}_{\text{peak}}$ and average read repair rate \mathcal{R}_{avg} needed to achieve a large MTDDL for small code systems can be substantial. Modifications of standard erasure codes have been designed for storage systems to reduce these rates, e.g., *local reconstruction codes* (LR codes) [9], [19], and *regenerating codes* (RG codes) [6], [7]. Some versions of LR codes have been used in deployments, e.g., by Microsoft Azure [10]. The encoding and decoding complexity of RS, LR, and RG codes grows non-linearly (typically quadratic or worse) as the values of (n, k, r) grow, which makes the use of such codes with large values of (n, k, r) less practical.

2.4. Liquid systems

We introduce *liquid systems* for reliably storing objects in a distributed set of storage nodes. Liquid systems use a combination of a *large code* with large values of (n, k, r) , *lazy repair*, and a *flow storage organization*. Because of the large code size the placement group concept is not of much importance: For the purposes of this paper we assume that one placement group is used for all objects, i.e., $n = M$ and a fragment is assigned to each node for each object.

The RaptorQ code [20], [13] is an example of an erasure code that is suitable for a liquid system. RaptorQ codes are fundamentally designed to support large values of (n, k, r) with very low encoding and decoding complexity and to have exceptional recovery properties. Furthermore, RaptorQ codes are *fountain codes*, which means that as many encoded symbols n as desired can be generated on the fly for a given value of k . The fountain code property provides flexibility when RaptorQ codes are used in applications, including liquid systems. The monograph [20] provides a detailed exposition of the design and analysis of RaptorQ codes.

The value of r for a liquid system is large, and thus *lazy repair* can be used, where lost fragments are repaired at a steady background rate using a reduced amount of bandwidth. The basic idea is that the read repair rate is controlled so that objects are typically repaired when a large fraction of r fragments are missing (ensuring high repair efficiency), but long before r fragments are missing (ensuring a large MTDDL), and thus the read repair rate is not immediately reactive to individual node failures. For a fixed node failure rate the read repair rate can in principle be fixed as described in Appendix A. When the node failure rate is unknown a priori, the algorithms described in Section 9 and Appendix A should be used to adapt the read repair rate to changes in conditions (such as node failure rate changes) to ensure high repair efficiency and a large MTDDL.

We show that a liquid system can be operated to enable flexible and essentially optimal combinations of storage durability, storage overhead, repair bandwidth usage, and access performance, exceeding the performance of small code systems.

Information theoretic limits on tradeoffs between storage overhead and the read repair rate are introduced and proved in [14]. The liquid systems for which prototypes and simulations are described

3 SYSTEM FAILURES

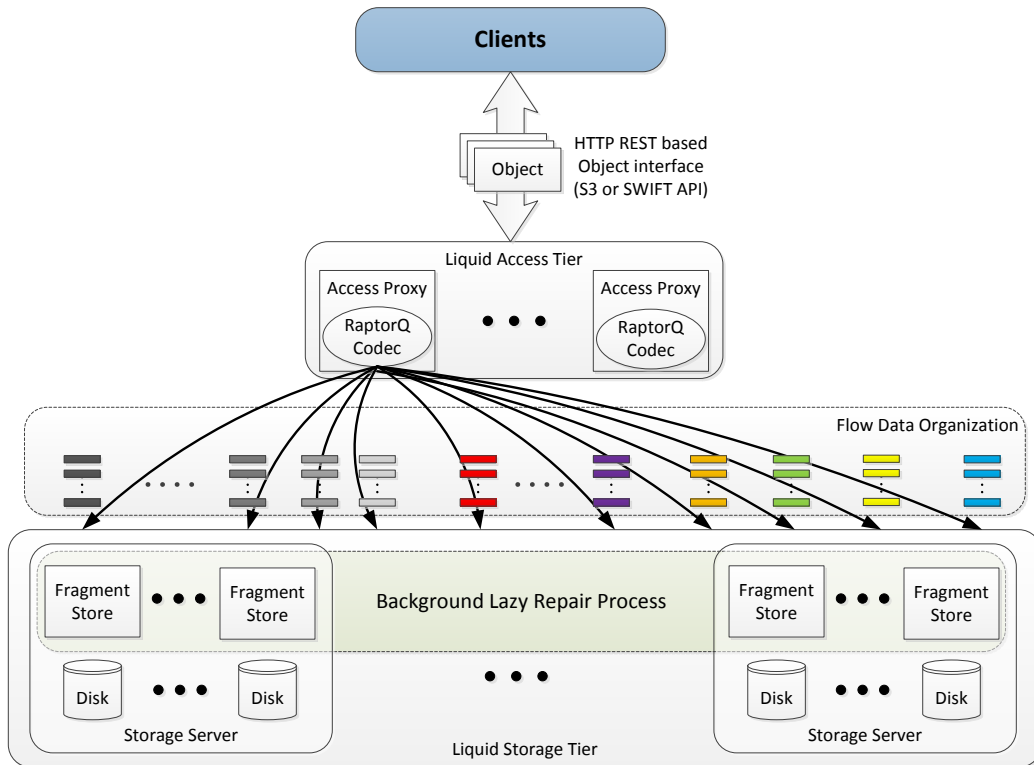


Fig. 1: A liquid system architecture.

in this paper are similar to the liquid systems described in Section 9.A of [14] that are within a factor of two of optimal. More advanced liquid systems similar to those described in Section 9.B of [14] that are asymptotically optimal would perform even better.

Fig. 1 shows a possible liquid system architecture. In this architecture, a standard interface, such as the S3 or SWIFT API, is used by clients to make requests to store and access objects. Object storage and access requests are distributed across access proxy servers within the liquid access tier. Access proxy servers use a RaptorQ codec to encode objects into fragments and to decode fragments into objects, based on the flow storage organization. Access proxy servers read and write fragments from and to the storage servers that store the fragments on disks within the liquid storage tier. Lazy repair operates as a background process within the liquid storage tier. Further details of this example liquid system architecture are discussed in Section 8.

3. SYSTEM FAILURES

Since distributed storage systems consist of a large number of components, failures are common and unavoidable. The use of commodity hardware can additionally increase the frequency of failures and outages, as can necessary maintenance and service operations such as software upgrades or system reconfiguration on the hardware or software level.

Different failure types can affect the system in various ways. For example, a hard drive can fail in its entirety, or individual sectors of it can become unreadable. Node failures can be permanent, e.g. if the failure results in the entire node needing to be replaced, or it can be transient, if it is caused issues such as network outages, reboots, software changes. Network failures can cause subsets of many nodes to be unreachable, e.g. a full rack or even a complete site.

Moreover, distributed storage systems are heterogeneous in nature, with components that are commonly being upgraded and replaced. Device models do not remain the same; for example, a replacement hard drive will often be of a different kind, with different and possibly unknown expected lifetime. Even within a device model, different manufacturing runs can result in variation in expected lifetimes. Further, upgrades may require a large set of machines to be replaced in a short time frame, leading to bursts of increased churn. Changes to the underlying technology, such as a transition from hard drives to SSDs, are also likely to profoundly impact the statistics and nature of storage failures.

For these reasons, the failure behavior of a distributed storage system is difficult to model. Results obtained under complicated realistic error models are difficult to interpret, and hard to generalize. The results in this paper are obtained under a fairly straightforward failure model which focuses on two types of failures: Node failures and sector failures.

3.1. Node failures

Although node failure processes can be difficult to accurately model, analysis based on a random node failure models can provide insight into the strengths and weaknesses of a practical system, and can provide a first order approximation to how a practical system will operate. The time till failure of an individual node is often modeled by a Poisson random variable with rate λ , in which case the average time till an individual node fails is $1/\lambda$. Nodes are assumed to fail independently. Since failed nodes are immediately replaced, the aggregate node failure process is Poisson with rate $\lambda \cdot M$. Our simulations use this model with $1/\lambda = 3$ years, and also use variants of this model where the node failure rate bursts to 3 times this rate ($1/\lambda = 1$ year) or 10 times this rate ($1/\lambda = 0.3$ years) for periods of time.

Nodes can become unresponsive for a period of time and then become responsive again (transient node failure), in which case the data they store is temporarily unavailable, or nodes can become and remain unresponsive (node failure), in which case the data they store is permanently unavailable. Transient node failures are an order of magnitude more common than node failures in practice, e.g., see [8].

Unresponsive node are generally detected by the system within a few seconds. However, it is usually unknown whether the unresponsiveness is due to a transient node failure or a node failure. Furthermore, for transient node failures, their duration is generally unknown, and can vary in an extremely large range: Some transient node failures have durations of just seconds when others can have durations of months, e.g., see [8].

Section 9 discusses more general and realistic node failure models and introduces repair strategies that automatically adjust to such models.

3.2. Sector failures

Sector failures are another type of data loss that occurs in practice and can cause operational issues. A sector failure occurs when a sector of data stored at a node degrades and is unreadable. Such data loss can only be detected when an attempt is made to read the sector from the node. (The data is typically stored with strong checksums, so that the corruption or data loss becomes evident when an attempt to read the data is made.) Although the sector failure data loss rate is fairly small, the overhead of detecting these types of failures can have a negative impact on the read repair rate and on the MTTDL.

Data *scrubbing* is often employed in practice to detect sector failures: An attempt is made to read through all data stored in the system at a regular frequency. As an example, [4] reports that Dropbox scrubs data at the frequency of once each two weeks, and reports that read traffic due to scrubbing can be greater than all other read data traffic combined. Most (if not all) providers of data storage systems scrub data, but do not report the frequency of scrubs. Obviously, scrubbing can have an impact on system performance.

The motivation for data scrubbing is that otherwise sector failures could remain undetected for long periods of time, and thus have a large negative impact on the MTTDL. This impact can be

4 REPAIR

substantial even if the data loss rate due to sector failures is relatively much smaller than the data loss rate due to node failures. The reason for this is that time scale for detecting and repairing sector failures is so much higher than for detecting and repairing node failures.

We do not employ a separate scrubbing process in our simulations, but instead rely on the repair process to scrub data: Since each node fails on average in $1/\lambda$ years, data on nodes is effectively scrubbed on average each $1/\lambda$ years. Sector failures do not significantly affect the MTDL of liquid systems, but can have a significant affect on the MTDL for small code systems.

We model the time to sector failure of an individual 4 KB sector of data on a node as a Poisson random variable with rate λ_{SF} , and thus the average time till an individual 4 KB sector of data on a node experiences a sector failure is $1/\lambda_{SF}$, where sector failures of different sectors are independent. Based on sector failure statistics derived from the data of [25], our simulations use this model with $1/\lambda_{SF} = 5 \cdot 10^8$ years. Thus, the data loss rate due to sector failures is more than 10^8 times smaller than the data loss rate due to node failures in our simulations.

4. REPAIR

Although using erasure coding is attractive in terms of reducing storage overheads and improving durability compared to replication, it can potentially require large amounts of bandwidth to repair for fragment loss due to node failures. Since the amount of required repair bandwidth can be high and can cause system-wide bottlenecks, this is sometimes referred to as a *repair bottleneck*. For example, [19] states the goal of their work as follows:

Our goal is to design more efficient coding schemes that would allow a large fraction of the data to be coded without facing this repair bottleneck. This would save petabytes of storage overheads and significantly reduce cluster costs.

The repair bottleneck sometimes has proven to be an obstacle in the transition from simple replication to more storage efficient and durable erasure coding in storage clusters. It simply refers to the fact that, when using a standard erasure code such as a RS code, the loss of a fragment due to a node failure requires the transfer of k fragments for repair, thus increasing k -fold the required bandwidth and I/O operations in comparison to replication. For example, [19] states the following when describing a deployed system with 3000 nodes:

... the repair traffic with the current configuration is estimated around 10-20% of the total average of 2 PB/day cluster network traffic. As discussed, (14,10,4) RS encoded blocks require approximately 10x more network repair overhead per bit compared to replicated blocks. We estimate that if 50% of the cluster was RS encoded, the repair network traffic would completely saturate the cluster network links.

The recognition of a repair bottleneck presented by the use of standard MDS codes, such as RS codes, has led to the search for codes that provide *locality*, namely where the repair requires the transfer of $\ell < k$ fragments for repair. The papers [9], [10], [19] introduce LR codes and provide tradeoffs between locality and other code parameters. We discuss some properties of small code systems based on LR codes in Section 5.2.

The following sub-sections provide the motivation and details for the repair strategies employed in our simulations.

4.1. Repair initiation timer

In practice, since the vast majority of unresponsive nodes are transient node failures, it is typical to wait for a *repair initiation time* T_{RT} after a node becomes unresponsive before initiating repair. For example, $T_{RT} = 30$ minutes for the small code system described in [10], and thus the policy reaction is not immediate but nearly immediate.

In our simulations repairers operate as follows. If a node is unresponsive for a period at most T_{RT} then the event is classified as a transient node failure and no repair is triggered. After being unresponsive for a period of T_{RT} time, a node failure is declared, the node is decommissioned, all

fragments stored on the node are declared to be missing (and will be eventually regenerated by the repairer). The decommissioned node is immediately replaced by a new node that initially stores no data (to recover the lost raw capacity).

Transient node failures may sometimes take longer than 30 minutes to resolve. Hence, to avoid unnecessary repair of transient node failures, it is desirable to set T_{rit} longer, such as 24 hours. A concern with setting T_{rit} large is that this can increase the risk of source data loss, i.e., it can significantly decrease the MTTDL.

Since a large value for T_{rit} has virtually no impact on the MTTDL for liquid systems, we can set T_{rit} to a large value and essentially eliminate the impact that transient node failures have on the read repair rate. Since a large value for T_{rit} has a large negative impact on the MTTDL for small code systems, T_{rit} must be set to a small value to have a reasonable MTTDL, in which case the read repair rate can be adversely affected if the unresponsive time period for transient node failure extends beyond T_{rit} and triggers unnecessary repair. We provide some simulation results including transient node failures in Section 5.4.

4.2. Repair strategies

In our simulations, a repairer maintains a repair queue of objects to repair. An object is added to the repair queue when at least one of the fragments of the object have been determined to be missing. Whenever there is at least one object in the repair queue the repair policy works to repair objects. When an object is repaired, at least k fragments are read to recover the object and then additional fragments are generated from the object and stored at nodes which currently do not have fragments for the object.

Objects are assigned to placement groups to determine which nodes store the fragments for the objects. Objects within a placement group are *prioritized* by the repair policy: objects with the *least* number of available fragments within a placement group are the next to be repaired within that placement group. Thus, objects assigned to the same placement group have a *nested object structure*: The set of available fragments for an object within a placement group is a subset of the set of available fragments for any other object within the placement group repaired more recently. Thus, objects within the same placement group are repaired in a round-robin order. Consecutively repaired objects can be from different placement groups if there are multiple placement groups.

4.3. Repair for small code system

In our simulations, small code systems (using either replication or small codes) use *reactive repair*, i.e., $\mathcal{R}_{\text{peak}}$ is set to a significantly higher value than the required average repair bandwidth, and thus repair occurs in a burst when a node failure occurs and the lost fragments from the node are repaired as quickly as possible. Typically only a small fraction of objects are in the repair queue at any point in time for a small code system using reactive repair.

We set the number of placement groups to $P = \frac{100 \cdot M}{n}$ to the recommended Ceph [22] value, and thus there are 100 placement groups assigned to each of the M nodes, with each node storing fragments for placement groups assigned to it. The repair policy works as follows:

- If there are at most 100 placement groups that have objects to repair then each such placement group repairs at a read repair rate of $\mathcal{R}_{\text{peak}}/100$.
- If there are more than 100 placement groups that have objects to repair then the 100 placement groups with objects with the least number of available fragments are repaired at a read repair rate of $\mathcal{R}_{\text{peak}}/100$.

This policy ensures that, in our simulations, the global peak read repair rate is at most $\mathcal{R}_{\text{peak}}$, that the global bandwidth $\mathcal{R}_{\text{peak}}$ is fully utilized in the typical case when one node fails and needs repair (which triggers 100 placement groups to have objects to repair), and that the maximum traffic from and to each node is not exorbitantly high. (Unlike the case for liquid systems, there are still significant concentrations of repair traffic in the network with this policy.)

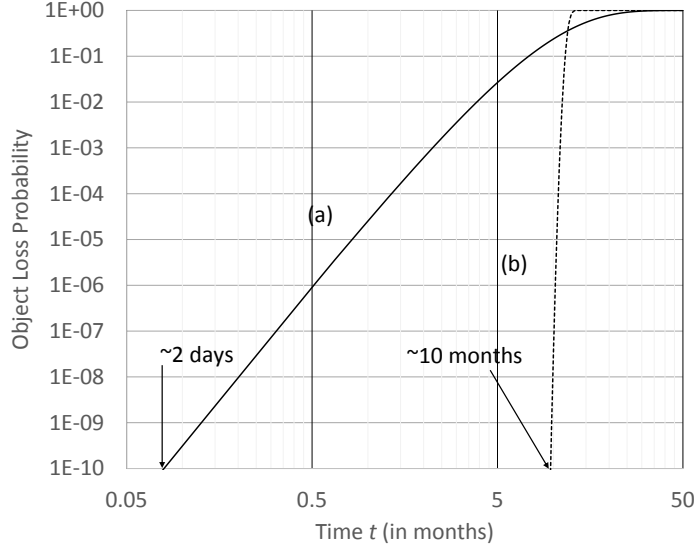


Fig. 2: Cumulative probability of object loss versus time for: (a) $(14, 10, 4)$ small code ; (b) $(3010, 2150, 860)$ large code.

The need for small code systems to employ reactive repair, i.e., use a high $\mathcal{R}_{\text{peak}}$ value, is a consequence of the need to achieve a large MTTDL. Small codes are quite vulnerable, as the loss of a small number of fragments can lead to object loss. For instance, with the $(14, 10, 4)$ RS code the loss of five or more fragments leads to object loss. Fig. 2 illustrates the drastically different time scales at which repair needs to operate for a small code system compared to a liquid system. In this example, nodes fail on average every three years and Fig. 2 shows the probability that $N \leq k - 1$ of n nodes are functioning at time t when starting with $N = n$ nodes functioning at time zero. As shown, a failed node needs to be repaired within two days to avoid a 10^{-10} probability object loss event for a $(14, 10, 4)$ RS code, whereas a failed node needs to be repaired within 10 months for a $(3010, 2150, 860)$ large code. Note that the storage overhead $\beta = r/n = 0.286$ is the same in these two systems. A similar observation was made in [21].

Exacerbating this for a small code system is that there are many placement groups and hence many combinations of small number of node failures that can lead to object loss, thus requiring a high $\mathcal{R}_{\text{peak}}$ value to ensure a target MTTDL. For example, $P = \frac{100 \cdot M}{n} = 21500$ for a $(14, 10, 4)$ small code in a $M = 3010$ node system, and thus an object from the system is lost if any object from any of the 21500 placement groups is not repaired before more than four nodes storing fragments for the object fail. Furthermore, a small code system using many placement groups is difficult to analyze, making it difficult to determine the value of $\mathcal{R}_{\text{peak}}$ that ensures a target MTTDL.

Although small codes outperform replication in both storage overhead and reliability metrics, they require a high $\mathcal{R}_{\text{peak}}$ value to achieve a reasonable MTTDL, which leads to the *repair bottleneck* described earlier. For example, $\mathcal{R}_{\text{peak}}$ used by small code system can be controlled to be a fraction of the total network bandwidth available in the cluster, e.g., 10% of the total available bandwidth is mentioned in [10] and 10% – 20% was estimated in [19] to protect a small fraction of the object data with a small code. However, using these fractions of available bandwidth for repair may or may not achieve a target MTTDL. As we show in Section 5 via simulation, small code systems require a large amount of available bandwidth for bursts of time in order to achieve a reasonable target

MTTDL, and in many cases cannot achieve a reasonable MTTDL due to bandwidth limitations, or due to operating with a large T_{RTT} value.

4.4. Repair for fixed rate liquid system

The simplest repair policy for a liquid system is to use a fixed read repair rate: the value of the peak read repair rate $\mathcal{R}_{\text{peak}}$ is set and the read repair rate is $\mathcal{R}_{\text{peak}}$ whenever there are objects in the repair queue to be processed. A liquid system employs a repairer that can be described as *lazy*: The value of $\mathcal{R}_{\text{peak}}$ is set low enough that the repair queue contains almost all objects all the time. Overall, a liquid system using lazy repair uses a substantially lower average repair bandwidth and a dramatically lower peak repair bandwidth than a small code system using reactive repair.

There are two primary reasons that a liquid system can use lazy repair: (1) the usage of a large code, i.e., as shown in Fig. 2, the (3010, 2150, 860) large code has substantially more time to repair than the (14, 10, 4) small code to achieve the same MTTDL; (2) the nested object structure.

A liquid system can store a fragment of each object at each node, and thus uses a single placement group for all objects. For example, using a (3010, 2150, 860) large code, an object is lost if the system loses more than 860 nodes storing fragments for the object before the object is repaired. Since there is only one placement group the nested object structure implies that if the object with the fewest number of available fragments can be recovered then all objects can be recovered. This makes it simpler to determine an $\mathcal{R}_{\text{peak}}$ value that achieves a target MTTDL for a fixed node failure rate. The basic intuition is that the repair policy should cycle through and repair the objects fast enough to ensure that the number of nodes that fail between repairs of the same object is at most r . Note that $\mathcal{D}_{\text{src}}/\mathcal{R}_{\text{peak}}$ is the time between repairs of the same object when the aggregate size of all objects is \mathcal{D}_{src} and the peak global read repair bandwidth is set to $\mathcal{R}_{\text{peak}}$. Thus, the value of $\mathcal{R}_{\text{peak}}$ should be set so that the probability that there are more than r node failures in time $\mathcal{D}_{\text{src}}/\mathcal{R}_{\text{peak}}$ is extremely tiny. Eq. (3) from Appendix A provides methodology for determining a value of $\mathcal{R}_{\text{peak}}$ that achieves a target MTTDL.

Unlike small code systems, wherein T_{RTT} must be relatively small, a liquid systems can use large T_{RTT} values, which has the benefit of practically eliminating unnecessary repair due to transient failures.

5. REPAIR SIMULATION RESULTS

The storage capacity of each node is $S = 1$ PB in all simulations. The system is fully loaded with source data, i.e., $\mathcal{D}_{\text{src}} = S \cdot M \cdot (1 - \beta)$, where β is the storage overhead expressed as a fraction.

There are node failures in all simulations and node lifetimes are modeled as independent exponentially distributed random variables with parameter λ . Unless otherwise specified, we set $1/\lambda$ to 3 years. When sector failures are also included, sector lifetimes are also exponentially distributed with $1/\lambda_{\text{SF}}$ fixed to $5 \cdot 10^8$ years.

At the beginning of each simulation the system is in a perfect state of repair, i.e., all n fragments are available for all objects. For each simulation, the peak global read repair bandwidth is limited to a specified $\mathcal{R}_{\text{peak}}$ value. The MTTDL is calculated as the number of years simulated divided by one more than the observed number of times there was an object loss event during the simulation. The simulation is reinitialized to a perfect repair state and the simulation is continued when there is an object loss event. To achieve accurate estimates of the MTTDL, each simulation runs for a maximum number of years, or until 200 object loss events have occurred.

Fig. 3 presents a comparison of repair traffic for a liquid system and a small code system. In this simulation, lazy repair for the liquid system uses a steady read repair rate of 704 Gbps to achieve an MTTDL of 10^8 years, whereas the reactive repair for the small code system uses read repair rate bursts of 6.4 Tbps to achieve an MTTDL slightly smaller than 10^7 years. For the small code system, a read repair rate burst starts each time a node fails and lasts for approximately 4 hours. One of the 3010 nodes fails on average approximately every 8 hours.

In Fig. 4, we plot $\mathcal{R}_{\text{peak}}$ against achieved MTTDL for a cluster of $M = 402$ nodes. The liquid systems use fixed read repair rate equal to the indicated $\mathcal{R}_{\text{peak}}$ and two cases of $\beta = 33.3\%$ and

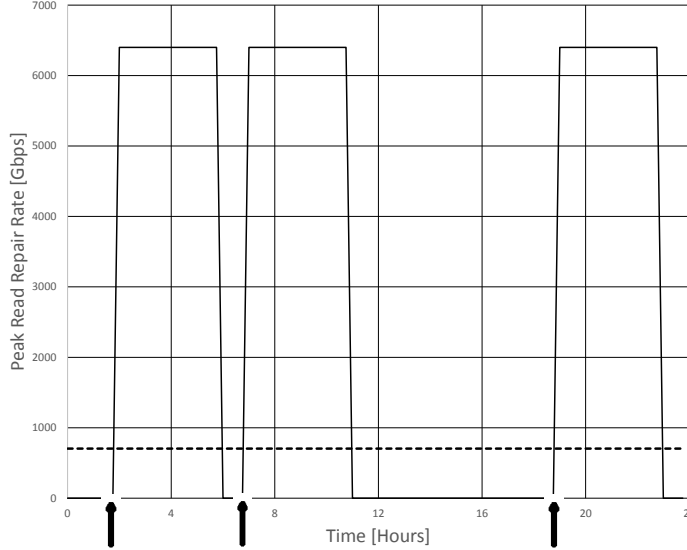


Fig. 3: Read repair rate comparison over a 24 hour period of reactive repair using a (14, 10, 4) small code (solid line) and lazy repair using a (3010, 2150, 860) large code (dash line). Arrows indicate times of node failures.

$\beta = 16.7\%$ are shown. The curves depict MTTDL bounds as calculated according to Eq. (3) of Appendix A. As can be seen, the bounds and the simulation results agree quite closely. Fig. 4 also shows simulation results for a small code system with storage overhead 33.3%, which illustrates a striking difference in performance between small code systems and liquid systems.

5.1. Fixed node failure rate with 402 nodes

Fig. 5 plots the simulation results for $\mathcal{R}_{\text{peak}}$ against MTTDL for a wide variety of $M = 402$ node systems. Each simulation was run for 10^9 years or until 200 object loss events occurred. The number of nodes, the node failure rate, and the small codes used in these simulations are similar to [10]. The square icons correspond to 16.7% storage overhead and the circle icons correspond to 33.3% storage overhead. For the small code systems, a (18, 15, 3) small code is used for 16.7% storage overhead, and a (9, 6, 3) small code is used for 33.3% storage overhead. For the liquid systems, a (402, 335, 67) large code is used for 16.7% storage overhead, and a (402, 268, 134) large code is used for 33.3% storage overhead.

The unshaded icons correspond to systems where there are only node failures, whereas the shaded icons correspond to systems where there are both node failures and sector failures.

The small code systems labeled as “SC, 30min” and “SC, 24hr” use the read repair rate strategy described in Section 4.3 based on the indicated $\mathcal{R}_{\text{peak}}$ value, with T_{RIT} set to 30 minutes and 24 hours respectively. The triplication systems labeled as “Trip, 30min” and “Trip, 24hr” use the read repair rate strategy described in Section 4.3 based on the indicated $\mathcal{R}_{\text{peak}}$ value with T_{RIT} is set to 30 minutes and 24 hours respectively.

The liquid systems labeled as “LiqF, 24hr” use a fixed read repair rate set to the shown $\mathcal{R}_{\text{peak}}$, which was calculated according to Eq. (3) from Appendix A with a target MTTDL of 10^7 years. The value of T_{RIT} is set to 24 hours. As would be expected since Eq. (3) is a lower bound, the observed MTTDL in the simulations is somewhat above 10^7 years.

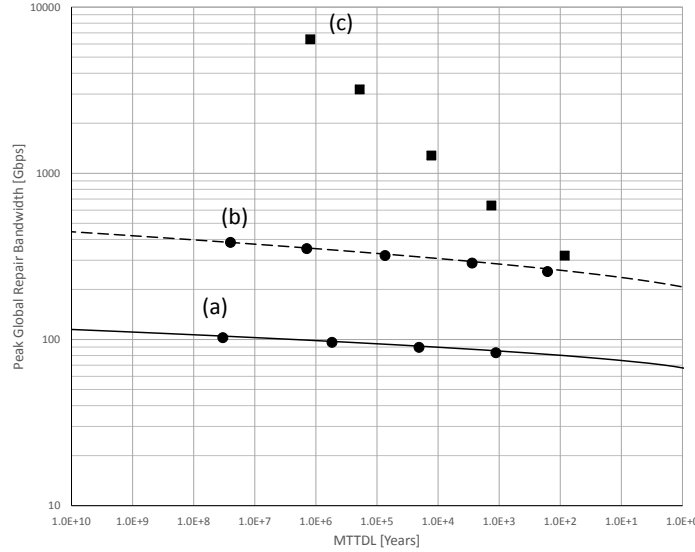


Fig. 4: Plot of MTTDL bound vs. $\mathcal{R}_{\text{peak}}$ for liquid systems with: (a) (402, 268, 134) code (33.3% storage overhead); (b) (402, 335, 67) code (16.7% storage overhead). Round markers on curve are simulation results. Square markers (c) are simulation results for small code system with (9, 6, 3) code (33.3% storage overhead).

The liquid systems labeled as “LiqR, 24hr” have T_{RTT} set to 24 hours and use a version of the regulator algorithm described in Section 9 to continually and dynamically recalculate the read repair rate according to current conditions. The regulator used (per object) node failure arrival rate estimates based on the average of the last $\frac{7}{6} \cdot r - F$ node failure inter-arrival times where F denotes the number of missing fragments for an object at the time it forms the estimate. The maximum read repair rate was limited to the shown $\mathcal{R}_{\text{peak}}$ value, which is about three times the average rate. In these runs there were no object loss events in the 10^9 years of simulation when there are both node failures and sector failures. This is not unexpected, since the estimate (using the MTTDL lower bound techniques described in Section 9 but with a slightly different form of node failure rate estimation) for the actual MTTDL is above 10^{25} years for the 33.3% storage overhead case and above $2.5 \cdot 10^{12}$ years for the 16.7% storage overhead case. These estimates are also known to be conservative when compared to simulation data. The same estimate shows that with this regulator an MTTDL of 10^9 years would be attained with the smaller overhead of $r = 115$ instead of $r = 134$. As shown in Table II, the average read repair rate \mathcal{R}_{avg} is around $1/3$ of $\mathcal{R}_{\text{peak}}$, and the 99% read repair rate $R_{99\%}$ is around $1/2$ of $\mathcal{R}_{\text{peak}}$. The large MTTDL achieved by the regulator arises from its ability to raise the repair rate when needed. Other versions of the regulator in which the peak read repair rate is not so limited can achieve much larger MTTDL. For example, a version with a read repair rate limit set to be about five times the average rate for the given node failure rate estimate achieves an MTTDL of around 10^{39} years for the 33.3% storage overhead case.

As can be seen, the value of $\mathcal{R}_{\text{peak}}$ required for liquid systems is significantly smaller than that required for small code systems and for triplication. Furthermore, although regulated rate liquid system can provide an MTTDL of greater than 10^{25} years even when T_{RTT} is set to 24 hours and there are sector failures in addition to node failures, the small code systems and triplication do not provide as good an MTTDL even when T_{RTT} is set to 30 minutes and there are no sector failures, and provide a poor MTTDL when T_{RTT} is set to 24 hours or when there are sector failures. The MTTDL

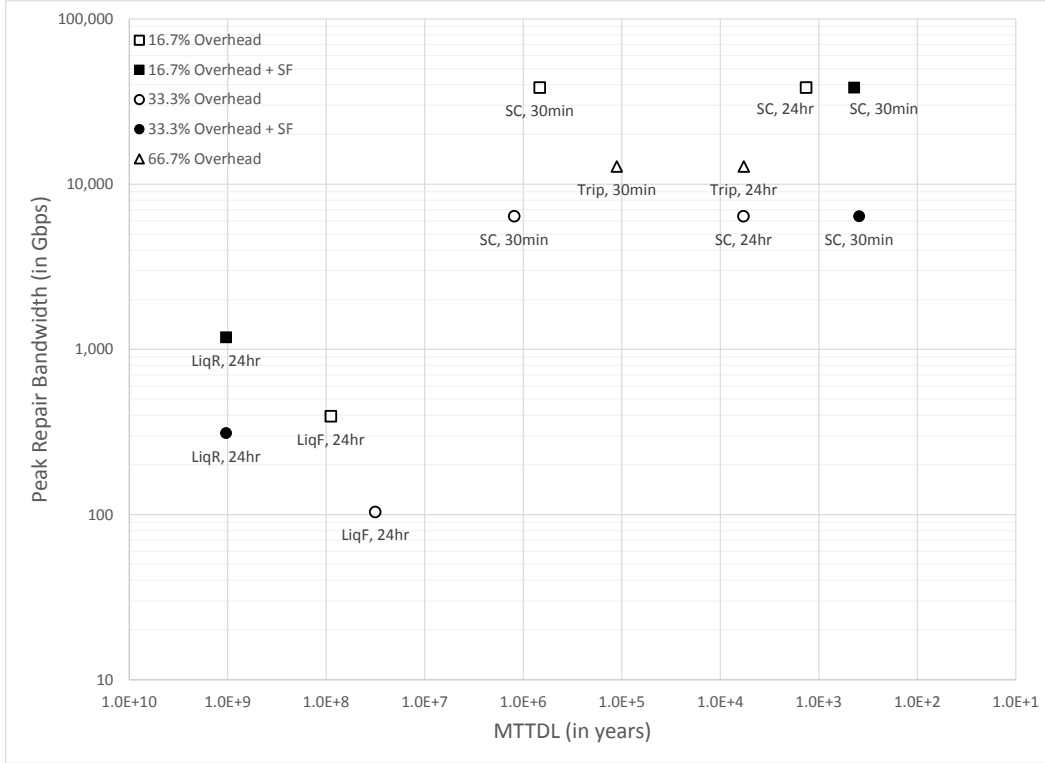


Fig. 5: Peak read repair rate versus MTTDL for a $M = 402$ node system for fixed $1/\lambda$.

would be further degraded for small code systems if both T_{RT} were set to 24 hours and there were sector failures.

5.2. Fixed node failure rate with 3010 nodes

Fig. 6 shows detailed simulation results, plotting the value of $\mathcal{R}_{\text{peak}}$ and the resulting value of MTTDL for various $M = 3010$ node systems. Each simulation was run for 10^8 years or until 200 object loss events occurred. The number of nodes, the node failure rate, and the small code systems used in these simulations are similar to [19]. The square icons correspond to 14.3% storage overhead and the circle icons correspond to 28.6% storage overhead. For the small code systems, a $(24, 20, 4)$ small code is used for 14.3% storage overhead, and a $(14, 10, 4)$ small code is used for 28.6% storage overhead. For the liquid systems, a $(3010, 2580, 430)$ large code is used for 14.3% storage overhead, and a $(3010, 2150, 860)$ large code is used for 28.6% storage overhead. The remaining parameters and terminology used in Fig. 6 are the same as in Fig. 5. There were no object loss events in the 10^8 years of simulation for any of the liquid systems shown in Fig. 6.

The systems “LiqR, 24hr” liquid systems used a regulated repair rate similar to that described for the 402 node system described above. The target repair efficiency was set at $\frac{2}{3}r$. In this case the average repair rate is higher than for the fixed rate case because the fixed rate case was set using the target MTTDL which, due the larger scale, admits more efficient repair than the 402 node system for the same target MTTDL. Estimates of MTTDL for the regulated case using techniques from Section 9 indicate that the MTTDL is greater than an amazing 10^{140} years for the 28.6% storage overhead case and 10^{66} years for the 14.3% storage overhead case even without the benefit of node failure rate estimation. As before, the regulated repair system could be run at substantially lower

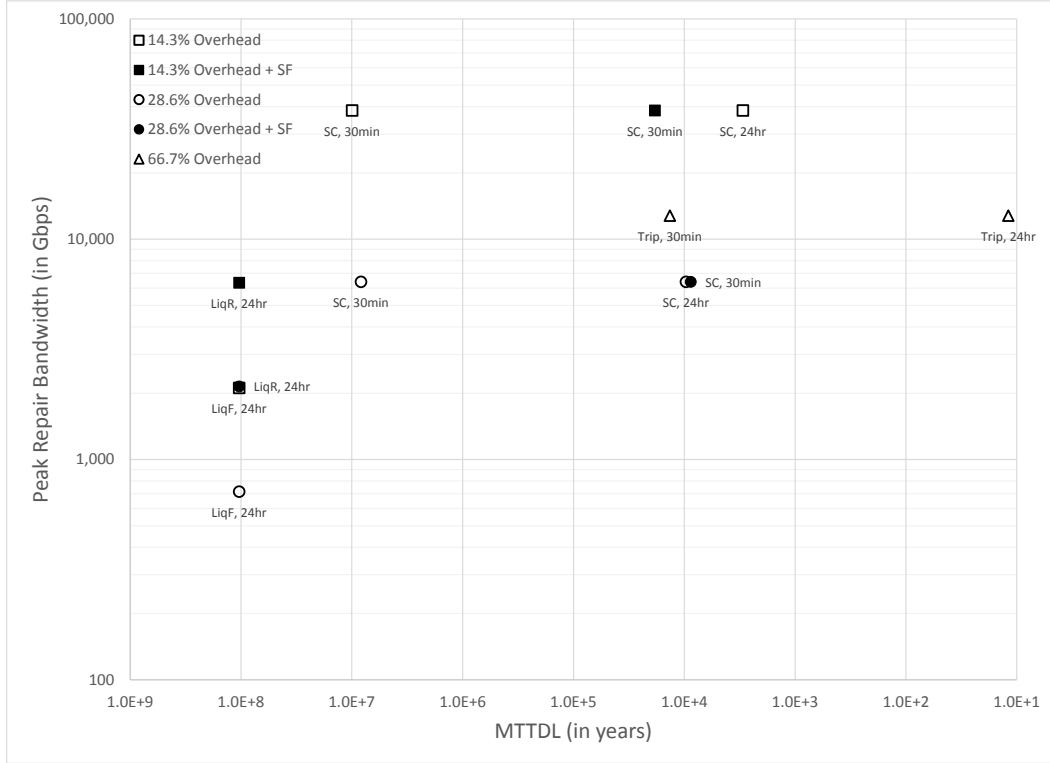


Fig. 6: Peak read repair rate versus MTTDL for a $M = 3010$ node system for fixed $1/\lambda$.

overhead while still maintaining a large MTTDL. For the “LiqR, 24hr” liquid systems, the average read repair rate \mathcal{R}_{avg} is around $1/3$ of $\mathcal{R}_{\text{peak}}$, and the 99% read repair rate $R_{99\%}$ is around $1/2$ of $\mathcal{R}_{\text{peak}}$.

The conclusions drawn from comparing small code systems to liquid systems shown in Fig. 5 are also valid for Fig. 6 but to an even greater degree. The benefits of the liquid system approach increase as the system size increases.

The paper [19] describes a *repair bottleneck* for a small code system, hereafter referred to as the Facebook system, which has around 3000 nodes and storage capacity per node of approximately $S = 15$ TB. The Facebook system uses a $(14, 10, 4)$ RS code to protect 8% of the source data and triplication to protect the remaining 92% of the source data. The amount of repair traffic estimated in [19] for the Facebook system is around 10% to 20% of the total average of 2 PB/day of cluster network traffic. The paper [19] projects that the network would be completely saturated with repair traffic if even 50% of the source data in the Facebook system were protected by the RS code.

Taking into account the mix of replicated and RS protected source data and the mixed repair traffic cost for replicated and RS protected source data, there are around 6.5 to 13 node failures per day in the Facebook system. An average of 6.5 to 13 node failures per day for a 3000 node system implies that each individual node on average fails in around 115 to 230 days, which is around 2.5 to 5 times the node failure rate of $1/\lambda = 3$ years we use in most of our simulations. If all source data were protected by a $(14, 10, 4)$ RS code in the Facebook system then the repair traffic per day would average around 1 to 2 PB. However, the node failure per day statistics in [19] have wide variations, which supports the conclusions in [19] that protecting most of the source data with the RS code in the Facebook system will saturate network capacity.

Note that 1 to 2 PB/day of repair traffic for the Facebook system implies an average read repair rate of around 90 to 180 Gbps. If the storage capacity per node were $S = 1$ PB instead of $S = 15$

TB then the Facebook system average read repair rate \mathcal{R}_{avg} would be approximately 6 to 12 Tbps, which is around 2.5 to 5 times the average read repair rate \mathcal{R}_{avg} for the very similar “SC, 30min” small code system shown in Fig. 6. This relative difference in the value of \mathcal{R}_{avg} makes sense, since the Facebook system node failure rate is around 2.5 to 5 times larger than for the “SC, 30min” small code system.

The “SC, 30min” small code system shown in Fig. 6 achieves an MTTDL of just under 10^7 years using a peak read repair rate $\mathcal{R}_{\text{peak}} = 6.4$ Tbps and using $T_{\text{RIT}} = 30$ minutes. Assume the node failure rate for the Facebook system is 5 times the node failure rate for the “SC, 30min” small code system shown in Fig. 6, and suppose we set $T_{\text{RIT}} = 30/5 = 6$ minutes for the Facebook system. The scaling observations in Section 5.5 show that this Facebook system achieves a MTTDL of just under $10^7/5 = 2 \cdot 10^6$ years when $\mathcal{R}_{\text{peak}} = 6.4 \cdot 5 = 32$ Tbps.

The average node failure rates reported in [19] for the Facebook system are considerably higher than the node failure rates we use in our simulations. One possible reason is that the Facebook system uses a small T_{RIT} value, e.g., $T_{\text{RIT}} = 10$ minutes, after which unresponsive nodes trigger repair. This can cause unnecessary repair of fragments on unresponsive nodes that would have recovered if T_{RIT} were larger.

It is desirable to set T_{RIT} as large as possible to avoid unnecessary repair traffic, but only if a reasonable MTTDL can be achieved. The results in Fig. 6 indicate that a small code system cannot offer a reasonable MTTDL with $T_{\text{RIT}} = 24$ hours. On the other hand, the results from Section 5.4 indicate that the average read repair rate increases significantly (indicating a significant amount of unnecessary repair) if a smaller value of T_{RIT} is used for a small code system when there are transient node failures. In contrast, a regulated liquid system can provide a very large MTTDL operating with $T_{\text{RIT}} = 24$ hours and with sector failures, thus avoiding misclassification of transient node failures and unnecessary repair traffic.

Results for small code systems using LR codes [9], [10], [19] can be deduced from results for small code systems using RS codes. A typical $(14, 10, 2, 2)$ LR code (similar to the one described in [10] which defines a $(16, 12, 2, 2)$ LR code) partitions 10 source fragments into two groups of five, generates one local repair fragment per group and two global repair fragments for a total of 14 fragments, and thus the storage overhead is the same as for a $(14, 10, 4)$ RS code. At least five fragments are read to repair a lost fragment using the LR code, whereas 10 fragments are read to repair a lost fragment using the RS code. However, there are fragment loss patterns the LR code does not protect against that the RS code does. Thus, the LR code operating at half the $\mathcal{R}_{\text{peak}}$ value used for the RS code achieves a MTTDL that is lower than the MTTDL for the RS code.

Similarly, results for [19] which use a $(16, 10, 4, 2)$ LR code can be compared against the $(14, 10, 4)$ RS code that we simulate. The LR code operating at half the $\mathcal{R}_{\text{peak}}$ value used for the RS code achieves an MTTDL that may be as large as the MTTDL for the RS code. However, the LR code storage overhead is $\beta = 0.375$, which is higher than the RS code storage overhead of $\beta = 0.286$.

5.3. Varying node failure rate

In this subsection we demonstrate the ability of the regulated repair system to respond to bursty node failure processes. The node failures in all simulations described in this subsection are generated by a time varying periodic Poisson process repeating the following pattern over each ten year period of time: the node failure rate is set to $1/\lambda = 3$ years for the first nine years of each ten year period, and then $1/\lambda = 1$ year for the last year of each ten year period. Thus, in each ten year period, the node failure rate is the same as it was in the previous subsections for the first nine years, followed by a node failure rate that is three times higher for the last year.

Except for using variable λ for node failures, Fig. 7 is similar to Fig. 5. The small code systems labeled as “SC, 30min” in Fig. 7 use the read repair rate strategy described in Section 4.3 based on the shown $\mathcal{R}_{\text{peak}}$ value, and T_{RIT} is set to 30 minutes. Thus, even when the node failure rate is lower during the first nine years of each period, the small code system still sets read repair rate to $\mathcal{R}_{\text{peak}}$ whenever repair is needed.

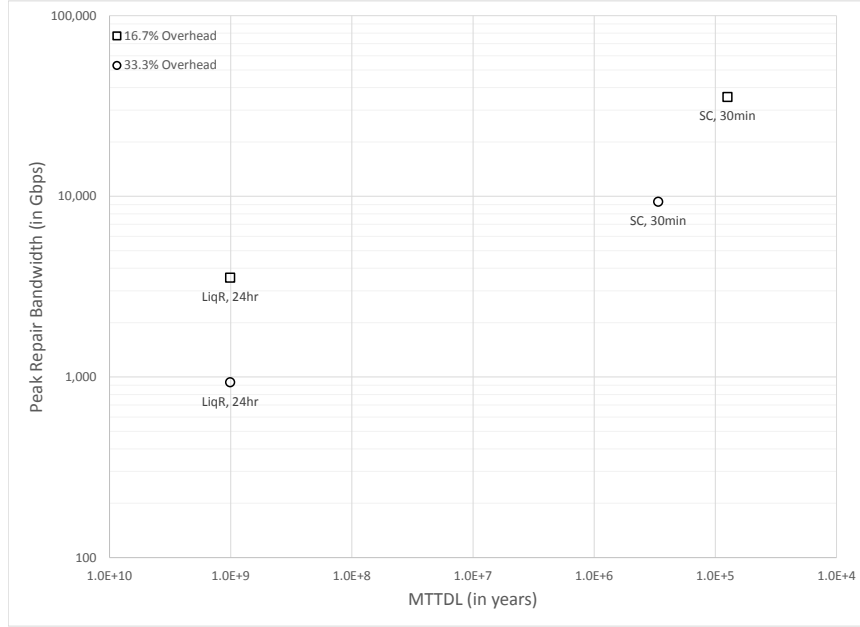


Fig. 7: Peak read repair rate versus MTTDL for a $M = 402$ node system for variable $1/\lambda$.

The liquid systems labeled as “LiqR, 24hr” in Fig. 7 use the regulator algorithm described in Section 9 to continually and dynamically recalculate the read repair rate according to current conditions, allowing a maximum read repair rate up to the shown $\mathcal{R}_{\text{peak}}$ value, and T_{RTT} is set to 24 hours. In these runs there were no object loss events in the 10^9 years of simulation with both node failures and sector failures. [[We may want to compute something here TBD?]]

Fig. 18 is an example trace of the read repair rate as function of time for the liquid system labeled “LiqR, 24hr” in Fig. 7, which shows how the the read repair rate automatically adjusts as the node failure rate varies over time. In these simulations, the average read repair rate \mathcal{R}_{avg} is around $1/9$ of $\mathcal{R}_{\text{peak}}$ during the first nine years of each period and around $1/3$ of $\mathcal{R}_{\text{peak}}$ during the last year of each period, and the 99% read repair rate $R_{99\%}$ is around $1/6$ of $\mathcal{R}_{\text{peak}}$ during the first nine years of each period and around $1/2$ of $\mathcal{R}_{\text{peak}}$ during the last year of each period. Thus, the regulator algorithm does a good job of matching the read repair rate to what is needed according to the current node failure rate.

Except for using variable λ for node failures, Fig. 8 is similar to Fig. 6. The conclusions comparing small code systems to liquid systems shown in Fig. 7 are also valid for Fig. 8.

5.4. Transient node failures

In this subsection we demonstrate the ability of the liquid systems to efficiently handle transient node failure processes.

Fig. 10 plots the simulation results for \mathcal{R}_{avg} against MTTDL for a number of $M = 402$ node systems. Like in previous subsections, the node lifetimes for node failures are modeled as independent exponentially distributed random variables with parameter λ and is set to $1/\lambda = 3$ years. The occurrence times of transient node failures are modeled as independent exponentially distributed random variables with parameter λ_{TF} and is set to $1/\lambda_{\text{TF}} = 0.33$ years. Thus, transient node failures occur at 9 times the rate at which node failures occur, consistent with [8]. The durations of transient node failures are modeled with log-logistic random variables, having a median of 60 seconds and shape parameter 1.1. These choices were made so as to mimic the distribution provided in [8]. Figure 9

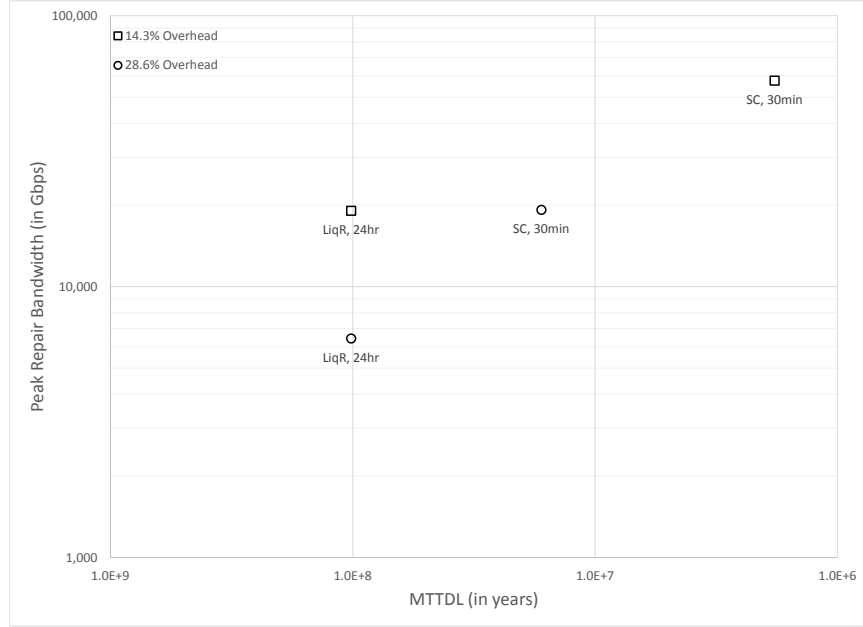


Fig. 8: Peak read repair rate versus MTTDL for a $M = 3010$ node system for variable $1/\lambda$.

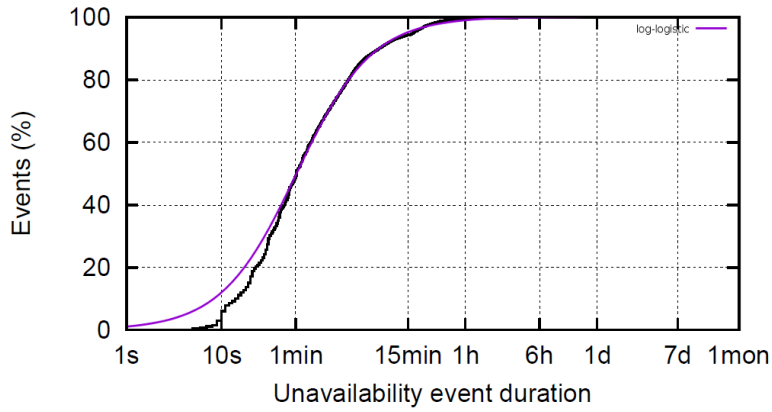


Fig. 9: Log-logistic transient failure duration model versus measured data

shows the graph from [8] with the log-logistic fitted curve overlaid. With this model, less than 10% of the transient node failures last for more than 15 minutes.

The unshaded markers mark simulations with just node failures and the shaded markers mark simulations with both node failures and transient node failures.

The $\mathcal{R}_{\text{peak}}$ values in all simulations is the same as $\mathcal{R}_{\text{peak}}$ value used correspondingly in Fig. 5. The small code systems labeled as “SC, 30min” and “SC, 24hr” in Fig. 10 use the read repair rate strategy described in Section 4.3, with $\mathcal{R}_{\text{peak}}$ value of 6400 Gbps and T_{RTT} is set to 30 minutes and 24 hours respectively. The liquid system labeled as “LiqR, 24hr” in Fig. 10 uses the regulator algorithm described in Section 9, with $\mathcal{R}_{\text{peak}}$ value of 311 Gbps and T_{RTT} is set to 24 hours.

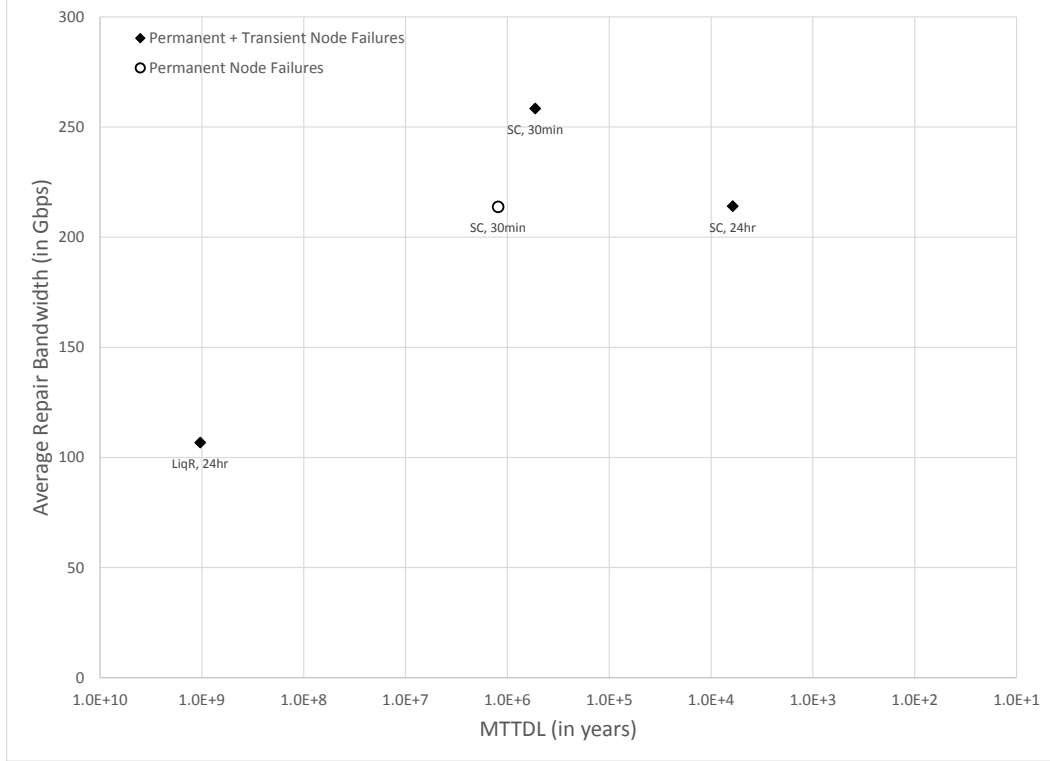


Fig. 10: Average read repair rate versus MTTDL for a $M = 402$ node system for transient $1/\lambda$.

As evident, there is no difference in the \mathcal{R}_{avg} for the liquid system between simulations with transient node failures and those without. No object loss events were observed for the liquid system in the 10^9 years of simulation with both node failures and transient node failures. The \mathcal{R}_{avg} for small code system labeled as “SC, 30min“ is however higher for the simulation with transient node failures and achieves an MTTDL that is less than half of what is achieved when there are no transient node failures. The \mathcal{R}_{avg} for small code system labeled as “SC, 24hr“ is the same between simulations with transient node failures and those without, but they fail to achieve a high enough MTTDL and do not provide adequate protection against object loss.

Thus, the liquid systems do a better job of handling transient node failures without requiring more effective average bandwidth than small code systems.

5.5. Scaling parameters

The simulation results shown in Section 5 are based on sets of parameters from published literature, and from current trends in deployment practices. As an example of a trend, a few years ago the storage capacity of a node was around $S = 16$ TB, whereas the storage capacity of a node in some systems is $S = 1$ PB or more, and thus the storage capacity per node has grown substantially.

Repair bandwidths scale linearly as a function of the storage capacity S per node (keeping other input parameters the same). For example, the values of $(\mathcal{R}_{\text{peak}}, \mathcal{R}_{\text{avg}})$ for a system with $S = 512$ TB can be obtained by scaling down by a factor of two the values of $(\mathcal{R}_{\text{peak}}, \mathcal{R}_{\text{avg}})$ from simulation results for $S = 1$ PB, whereas the MTTDL and other values remain unchanged.

Repair bandwidths and the MTTDL scale linearly as a function of concurrently scaling the node failure rate λ and the repair initiation timer T_{RIT} . For example, the values of $(\mathcal{R}_{\text{peak}}, \mathcal{R}_{\text{avg}})$ for a system with $1/\lambda = 9$ years and $T_{\text{RIT}} = 72$ hours can be obtained by scaling down by a factor of three the

6 ENCODING AND DECODING OBJECTS

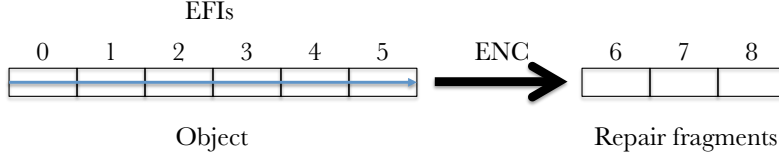


Fig. 11: Example of block storage organization. A data object is segmented into $k = 6$ source fragments and encoded by a small code to generate $r = 3$ repair fragments.

values of $(\mathcal{R}_{\text{peak}}, \mathcal{R}_{\text{avg}})$ from simulation results for $1/\lambda = 3$ years and $T_{\text{RTT}} = 24$ hours, whereas the MTDL can be obtained by scaling up by a factor of three, and other values remain unchanged.

6. ENCODING AND DECODING OBJECTS

In a *block storage organization* commonly used by small code systems, objects are partitioned into k successive source fragments and encoded to generate $r = n - k$ repair fragments (we assume the small code is MDS). In effect a fragment is a symbol of the small code, although the small code may be intrinsically defined over smaller symbols, e.g., bytes. A high level representation of a block storage organization for a simple example using a $(9, 6, 3)$ Reed-Solomon code is shown in Fig. 11, where the object runs from left to right across the source fragments as indicated by the blue line.

A relatively small chunk of contiguous data from the object that resides within a single source fragment can be accessed directly from the associated node if it is available. If, however, that node is unavailable, perhaps due to a transient node failure or node failure, then k corresponding chunks of equal size must be read from other nodes in the placement group and decoded to generate the desired chunk. This involves reading k times the size of the missing chunk and performing a decoding operation, which is referred to as a *degraded read* [19], [18]. Reading k chunks can incur further delays if any of the nodes storing the chunks are busy and thus non-responsive. This latter situation can be ameliorated if a number of chunks slightly larger than k are requested and decoding initiated as soon as the first k chunks have arrived, as described in [11].

A *flow storage organization*, used by liquid systems, operates in a stream fashion. An (n, k, r) large code is used with small symbols of size $Ssize$, resulting in a small source block of size $Bsize = k \cdot Ssize$. For example, with $k = 1024$ and symbols of size $Ssize = 64$ Bytes, a source block is of size $Bsize = 64$ KB. An object of size $Osize$ is segmented into

$$N = \frac{Osize}{Bsize}$$

source blocks, and the k source symbols of each such source block is independently erasure encoded into n symbols. For each $i = 0, \dots, n - 1$, the fragment with EFI i consists of the concatenation of the i^{th} symbol from each of the N consecutive source blocks of the object. An example of a flow storage organization showing a $(1536, 1024, 512)$ large code with $N = 24, 576$ is shown in Fig. 12, where each row corresponds to a source block and runs from left to right across the source symbols as indicated by the dashed blue lines, and the object is the concatenation of the rows from top to bottom.

A chunk of data consisting of a consecutive set of source blocks can be accessed as follows. For a fragment with EFI i , the consecutive portion of the fragment that corresponds to the i^{th} symbol from each of consecutive set of source blocks is read. When such portions from at least k fragments are received (each portion size a $1/k$ -fraction of the chunk size), the chunk of data can be recovered by decoding the consecutive set of source blocks in sequence. Thus, the amount of data that is read to access a chunk of data is equal to the size of the chunk of data. Portions from slightly more than k fragments can be read to reduce the latency due to nodes that are busy and thus non-responsive. Note that each access requires a decoding operation.

Let us contrast the two organizations with a specific example using an object of size $Osize = 1.5$ GB. Consider an access request for the 32 MB chunk of the object in positions 576 MB through

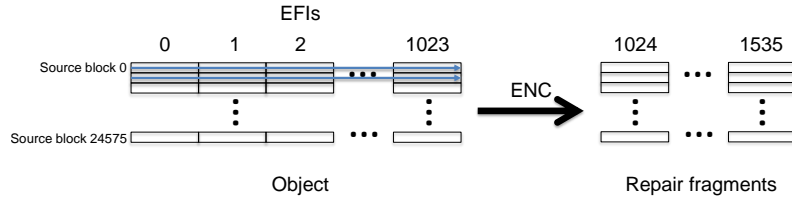


Fig. 12: Example of a flow storage organization. A data object is segmented into 24,576 source blocks, and each source block is segmented into $k = 1024$ source symbols and encoded by a large code to generate $r = 512$ repair symbols. The concatenation of the i^{th} symbol from each of the source blocks forms the fragment with EFI i . Thus the vertical columns labeled by the EFIs correspond to fragments.

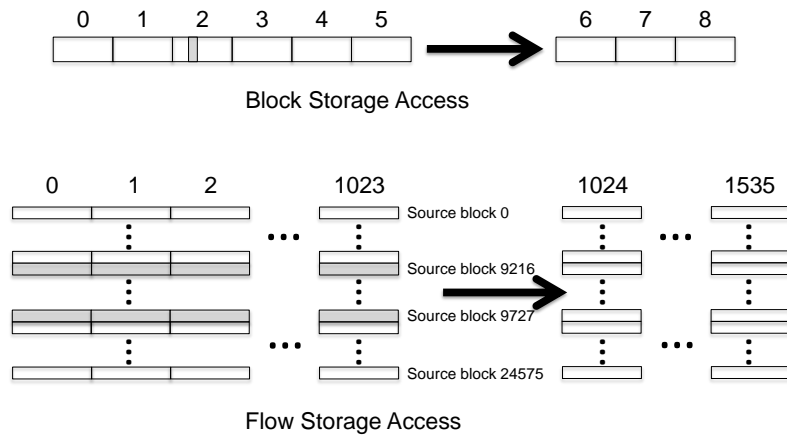


Fig. 13: Comparison of accessing a 32 MB chunk of an object for block storage organization and flow storage organization. The shaded region indicates the chunk of the object to be accessed.

608 MB, which is shown as the shaded region in Fig. 13 with respect to both a block storage organization and a flow storage organization.

With the block storage organization shown in Fig. 11, the 32 MB chunk is within the fragment with EFI = 2 as shown in Fig. 13. If the node that stores this fragment is available then the chunk can be read directly from that node. If that node is unavailable, however, then in order to recover the 32 MB chunk the corresponding 32 MB chunks must be retrieved from any six available nodes of the placement group and decoded. In this case a total of 192 MB of data is read and transferred through the network to recover the 32 MB chunk of data.

With the flow storage organization shown in Fig. 12, the chunk of size 32 MB corresponds to 512 consecutive source blocks from source block 9216 to block 9727 as shown in Fig. 13, which can be retrieved by reading the corresponding portion of the fragment from each of at least 1024 fragments. From the received portions, the chunk can be recovered by decoding the 512 source blocks. Thus, in a flow storage organization a total of 32 MB of data is read and transferred through the network to recover the original 32 MB chunk of data.

Naturally, the symbol size, source block size, and the granularity of the data to be accessed should be chosen to optimize the overall design, taking into account constraints such as whether the physical storage is hard disk or solid state.

7 PROTOTYPE IMPLEMENTATION

7. PROTOTYPE IMPLEMENTATION

Our team implemented a prototype of the liquid system. We used the prototype to understand and improve the operational characteristics of a liquid system implementation in a real world setting, and some of these learnings are described in Section 8. We used the prototype to compare access performance of liquid systems and small code systems. We used the prototype to cross-verify the lazy repair simulator that uses a fixed read repair rate as described in Section 4.4 to ensure that they both produce the same MTDL under the same conditions. We used the prototype to validate the basic behavior of a regulated read repair rate as described in Section 9.

The hardware we used for our prototype consists of a rack of 14 servers. The servers are connected via 10 Gbit full duplex ethernet links to a switch, and are equipped with Intel Xeon CPUs running at 2.8 GHz, with each CPU having 20 cores. Each server is equipped with an SSD drive of 600 GB capacity that is used for storage.

The liquid prototype system software is custom written in the Go programming language. It consists of four main modules:

- *Storage Node (SN)* software. The storage node software is a HTTP server. It is used to store fragments on the SSD drive. Since a liquid system would generally use many more than 14 storage nodes, we ran many instances of the storage node software on a small number of physical servers, thereby emulating systems with hundreds of nodes.
- An *Access Generator* creates random user requests for user data and dispatches them to accessors. It also collects the resulting access times from the accessors.
- *Accessors* take user data requests from the access generator and create corresponding requests for (full or partial) fragments from the storage nodes. They collect the fragment responses and measure the amount of time it takes until enough fragments have been received to recreate the user data. There are in general multiple accessors running in the system.
- A *Repair Process* is used to manage the repair queue and regenerate missing fragments of stored objects. The repair process can be configured to use a fixed read repair rate, or to use a regulated read repair rate as described in Section 9.

In addition, the team developed a number of utilities to test and exercise the system, e.g., utilities to cause node failures, and tools to collect data about the test runs.

7.1. Access performance setup

The users of the system are modeled by the access generator module. We set a target user access load ρ ($0 < \rho < 1$) on the system as follows. Let C be the aggregate capacity of the network links to the storage nodes, and let s be the size of the user data requests that will be issued. Then C/s requests per unit of time would use all available capacity. We set the mean time between successive user requests to $t = \frac{s}{\rho \cdot C}$, so that on average a fraction ρ of the capacity is used. The access generator module uses an exponential random variable to generate the timing of each successive user data request, where the mean of the random variable is t .

As an example, there is a 10 Gbps link to each of the six storage nodes in the setup of Fig. 14, and thus $C = 60$ Gbps. If user data requests are of size $s = 10$ MB = $80 \cdot 2^{20}$ bits, and the target load is $\rho = 0.8$ (80% of capacity), then the mean time between requests is set to

$$t = \frac{80 \cdot 2^{20}}{0.8 \cdot 60 \cdot 10^9} \approx 1.75 \text{ milliseconds.}$$

The access generator module round-robins the generated user data requests across the accessor modules. When an accessor module receives a generated user data request from the access generator module, the accessor module is responsible for making fragment requests to the storage nodes and collecting response payloads; it thus needs to be fairly high performance. Our implementation of an accessor module uses a custom HTTP stack, which pipelines requests and keeps connections statically alive over extended periods of time. Fragment requests are timed to avoid swamping the

switch with response data and this reduces the risk of packet loss significantly. The server network was also tuned, e.g., the TCP retransmission timer was adjusted for the use case of a local network. See [5] for a more detailed discussion of network issues in a cluster environment. These changes in aggregate resulted in a configuration that has low response latency and is able to saturate the network links fully.

7.2. Access performance tests

The prototype was used to compare access performance of liquid systems and small code systems. Our goal was to evaluate access performance using liquid system implementations, i.e., understand the performance impact of requesting small fragments (or portions of fragments) from a large number of storage servers. For simplicity we evaluated the network impact on access performance. Fragment data was generally in cache and thus disk access times were not part of the evaluation. We also did not evaluate the computational cost of decoding. Some of these excluded aspects are addressed separately in Section 8.

Figure 14 displays the setup that we use for testing access speed. Each physical server has a full duplex 10 GbE network link which is connected over a switch to all the other server blades. Testing shows that we are able to saturate all of the 10 Gbps links of the switch simultaneously. Eight of our 14 physical servers are running accessors, and one of the eight is running the access generator as well. The remaining six servers are running instances of the storage node software. This configuration allows us to test the system under a load that saturates the network capacity to the storage nodes: The aggregate network capacity between the switch and the accessors is 80 Gbps, which is more than the aggregate network capacity of 60 Gbps between the switch and the storage nodes, so the bottleneck is the 60 Gbps between the switch and the storage nodes.

During access performance tests there are no node failures, all n fragments for each object are available, and the repair process is disabled. For all tests, 67 storage node instances run on each storage server, which emulates a system of 402 nodes. We operate with a storage overhead of 33.3%, a (402, 268, 134) large code is used for the liquid system, and a (9, 6, 3) small code is used for the small code system. We tested with 10 MB and 100 MB user data request sizes.

We run the access prototype in three different configurations. The “Liq” configuration models liquid system access of user data. For a user data request of size s , the accessor requests $k + E$ fragment portions, each of size s/k , from a random subset of distinct storage nodes, and measures the time until the first k complete responses are received (any remaining up to E responses are discarded silently by the accessor). Each fragment portion is of size around 25.5 KB when $s = 10$ MB, and around 255 KB when $s = 100$ MB. We use $E = 30$, and thus the total size of requested fragment portions is around 10.75 MB when $s = 10$ MB, and around 107.5 MB when $s = 100$ MB. See [11] for an example of this access strategy.

The “SC” configuration models small code system normal access of user data, i.e., the requested user data is stored on a storage node that is currently available. For a user data request of size s , the accessor requests one fragment portion of size s from a randomly selected storage node, and measures the time until the complete response is received. The total size of the requested fragment portion is 10 MB when $s = 10$ MB, and 100 MB when $s = 100$ MB.

The “SC-Deg” configuration models small code system *degraded access* of user data, i.e., the requested user data is stored at a storage node that is currently unavailable, e.g., see [19], [18]. For a user data request of size s , the accessor requests k fragment portions of size s from a random subset of distinct storage nodes and measures the time until the first k complete responses are received. The total size of requested fragment portions is 60 MB when $s = 10$ MB, and 600 MB when $s = 100$ MB. Most user data is stored at available storage nodes, and thus most accesses are normal, not degraded, in operation of a small code system. Thus, we generate the desired load on the system with normal accesses to user data, and run degraded accesses at a rate that adds only a nominal aggregate load to the system, and only the times for the degraded accesses are measured by the accessors. Similar to the “Liq” configuration, we could have requested $k + 1$ fragment portions

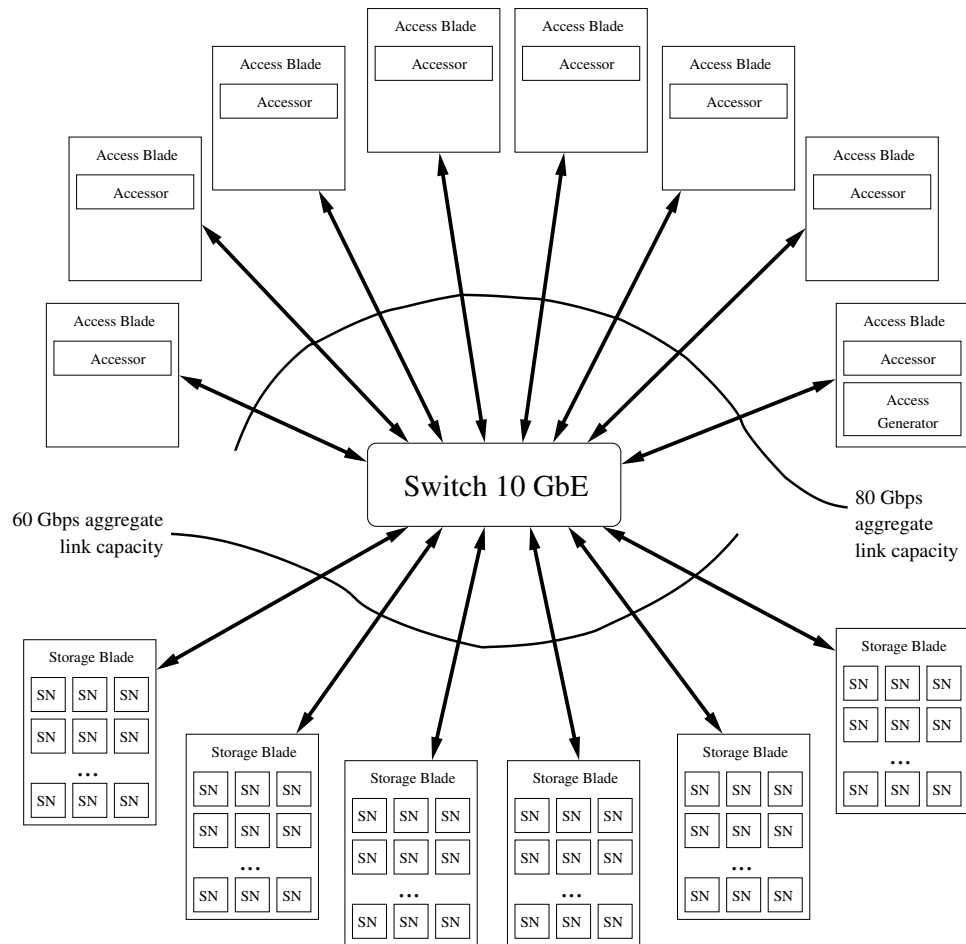


Fig. 14: Setup used for access performance testing

for the “SC-Deg” configuration to decrease the variation in access times, but at the expense of even more data transfer over the network.

Figures 15 and 16 depict access time results obtained with the above settings under different system loads. The average time for liquid system accesses is less than for small code system normal accesses in most cases, except under very light load when they are similar on average. Even more striking, the variation in time is substantially smaller for liquid systems than for small code systems under all loads. Liquid system accesses are far faster (and consume less network resources) than small code system degraded accesses.

Our interpretation of these results is that liquid systems spread load equally, whereas small code systems tend to stress individual storage nodes unevenly, leading to hotspots. With small code systems, requests to heavily loaded nodes result in slower response times. Hotspots occur in small code systems when appreciably loaded, even when the fragment requests are distributed uniformly at random. It should be expected that if the requests are not uniform, but some data is significantly more popular than other data, then response times for small code systems would be even more variable.

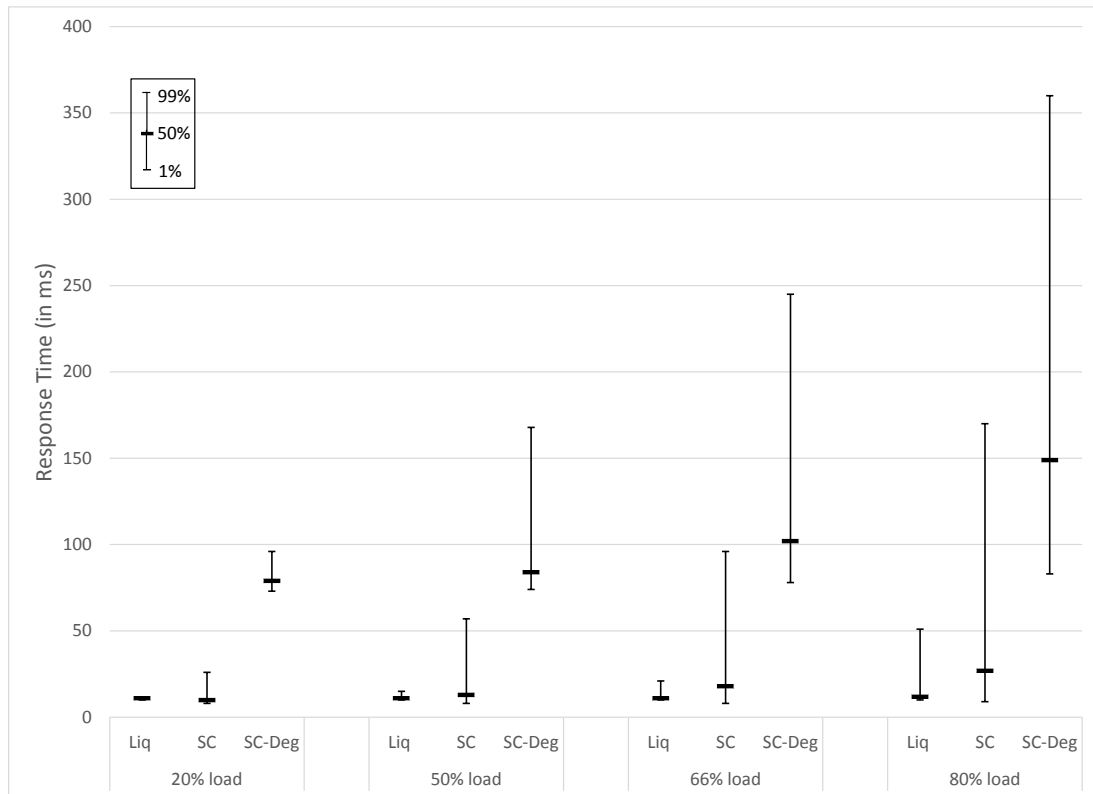


Fig. 15: Access Results for 10 MB object requests.

7.3. Repair tests and verification of the simulator

We used the prototype repair process to verify our liquid system repair simulator. We started from realistic sets of parameters, and then sped them up by a large factor (e.g., 1 million), in such a way that we could run the prototype repair process and observe object loss in a realistic amount of time. We ran the prototype repair process in such a configuration, where node failures were generated artificially by software. The resulting measured statistics, in particular the MTDL, were compared to an equivalent run of the liquid system repair simulator. This allowed us to verify that the liquid system repair simulator statistics matched those of the prototype repair process and also matched our analytical predictions.

8. IMPLEMENTATION CONSIDERATIONS

Designing a practical liquid system poses some challenges that are different from those for small code systems. In the following sections we describe some of the issues, and how they can be addressed.

8.1. Example Architecture

As discussed in Section 2.4, Fig. 1 shows an example of a liquid system architecture. This simplified figure omits components not directly related to source data storage, such as system management components or access control units.

Concurrent capacity to access and store objects scales with the number of access proxy servers, while source data storage capacity can be increased by adding more storage nodes to the system, and

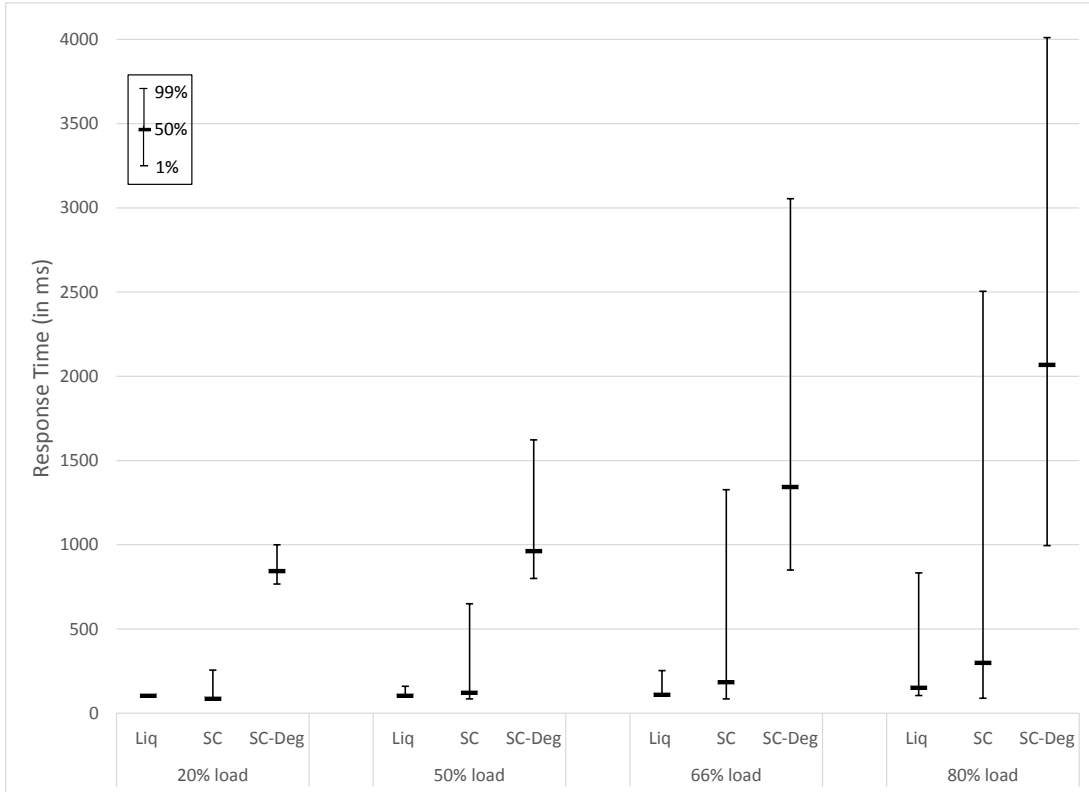


Fig. 16: Access Results for 100 MB object requests.

thus these capacities scale independently with this architecture. Storage nodes (storage servers) can be extremely simple, as their only function is to store and provide access to fragments, which allows them to be simple and inexpensive. The lazy repair process is independent of other components; thus reliability is handled by a dedicated system.

8.2. Metadata

The symbol size used by a liquid system with a flow storage organization is substantially smaller than the size used by a small code system with a block storage organization. For a small code system each symbol is a fragment, whereas for a liquid system there are many symbols per fragment and, as described in Section 6, the symbol to fragment mapping can be tracked implicitly. For both types of systems, each fragment can be stored at a storage node as a file.

The names and locations of each fragment are explicitly tracked for small code systems. There are many more fragments to track for a liquid system because of the use of a large code, and explicit tracking is less appealing. Instead, a mapping from the placement group to available nodes can be tracked, and this together with a unique name for each object can be used to implicitly derive the names and locations of fragments for all objects stored at the nodes for liquid systems.

8.3. Network usage

The network usage of a liquid system is fairly different from that of a small code system. As described in Section 7.2, liquid systems use network resources more smoothly than small code systems, but liquid systems must efficiently handle many more requests for smaller data sizes. For example, a small code system using TCP might open a new TCP connection to the storage node for

each user data request. For a liquid system, user data requests are split into many smaller fragment requests, and opening a new TCP connection for each fragment request is inefficient. Instead, each access proxy can maintain permanent TCP connections to each of the storage nodes, something that modern operating systems can do easily. It is also important to use protocols that allow pipeline data transfers over the connections, so that network optimizations such as jumbo packets can be effectively used. For example the HTTP/2 protocol [1] satisfies these properties.

Technologies such as RDMA [16] can be used to eliminate processing overhead for small packets of data. For example, received portions of fragments can be placed autonomously via RDMA in the memory of the access proxy, and the access proxy can decode or encode one or more source blocks of data concurrently.

8.4. Storage medium

A key strength of liquid systems is that user data can be accessed as a stream, i.e., there is very little startup delay until the first part of a user data request is available, and the entire user data request is available with minimal delay, e.g., see Section 7.2. Nevertheless, random requests for very small amounts of user data is more challenging for a liquid system.

The storage medium used at the storage nodes is an important consideration when deciding the large code parameters (n, k, r) . Hard disk drives (HDDs) efficiently support random reads of data blocks of size 500 KB or larger, whereas solid state drives (SSDs) efficiently support random reads of data blocks of size 4 KB or larger. We refer to these data block sizes as the *basic read size* of the storage medium. The parameters (n, k, r) should be chosen so that if s is the typical requested user data size then s/k is at least the basic read size, which ensures that a requested fragment portion size is at least the basic read size. In many cases this condition is satisfied, e.g., when the requested user data size is generally large, or when the storage medium is SSD.

In other cases, more advanced methods are required to ensure that a liquid system is efficient. For example, user data request performance can be improved by assigning different size user data to different types of objects: Smaller user data is assigned to objects that are stored using smaller values of k , and larger user data is assigned to objects that are stored using larger values of k , where in each case s/k is at least the basic read size if s is the requested user data size. In order to achieve a large MTTDL for all objects, the storage overhead $\frac{r}{n}$ is set larger for objects stored using smaller values of k than for objects stored using larger values of k . Overall, since the bulk of the source data is likely to be large user data, the overall storage overhead remains reasonable with this approach.

Only a small proportion of user data is frequently accessed in many use cases, e.g., in the use case described in [15], only about 10% of the user data is accessed frequently, and this user data is easy to identify based on the amount of time it has been stored in the storage system. Thus, a caching layer can be used to store and make such user data available without having to request the user data from the storage system. For a liquid system, it makes most sense to cache user data that is both small and frequently accessed.

Another strategy to minimize requests for small user data is to store related user data together in objects and coalesce accesses. Suppose for example, the storage system stores web content. To load a web page, dozens of small files of user data are typically loaded, such as HTML content and small images. Since the same content is accessed whenever such a web page is displayed, the system can store and access this user data together.

8.5. RaptorQ software performance

With small code systems, user data can usually be accessed directly without requiring erasure decoding. With liquid systems, each access of user data typically requires erasure decoding.

RaptorQ codes are the large codes of choice for liquid systems. The encoding and decoding complexity of RaptorQ codes is inherently linear by design [20], i.e., for a (n, k, r) RaptorQ code, the encoding complexity is linear in n , and the decoding complexity is linear in k , which makes the use of a RaptorQ code with large values of (n, k, r) practical.

Number of EFIs used to decode	1000	1001	1002	1003	1004
Fraction of failed decodings	$4.9 \cdot 10^{-3}$	$2.4 \cdot 10^{-5}$	$1.3 \cdot 10^{-7}$	$1.1 \cdot 10^{-9}$	0

Table I: Decode failures in 10^{11} decodings of a (1200, 1000, 200) RaptorQ code

Software implementations of RaptorQ codes have been deployed in a number of applications, including real-time applications that require low encoding and decoding complexity. Encode and decode speeds of over 10 Gbps are achieved using a software implementation of RaptorQ running on one core of a 20 core Intel Xeon CPU ES-2680 server running at 2.8 GHz, with code parameters $(n, k, r) = (1500, 1000, 500)$ and 256 byte symbols. Aggregate encode and decode speeds of around 150 Gbps are achieved using all 20 cores of the same server, with the same code parameters and symbol size.

These numbers suggest that RaptorQ decoding would add a fractional extra cost to the data processing pipeline: It is expected that other transformations are performed when accessing user data, such as for example checksumming, decompression and decryption. Of those MD5 checksums can be computed on the same machine at a rate of at most about 4 Gbps. As for decompression, Google’s run-time performance oriented Snappy codec achieves reportedly about 4 Gbps on a Core i7 processor [23], suggesting that RaptorQ coding is not a bottleneck.

8.6. RaptorQ hardware performance

Hardware based implementations have the potential to further cut down on power usage by the RaptorQ codes. Based on a hardware RaptorQ code chip design by Qualcomm for a different application, our estimate is that the power usage in an ASIC implementation would be about 3 Watts per 100 Gbps encoding or decoding speed.

To put this into perspective, we compare this to the power use of an SSD to access data. For example, the Samsung SM863a 2 TB enterprise SSD drive [24], achieves a data transfer rate of up to about 4.1 Gbps and uses about 2.5 W of power when reading. This translates to a power use of about 61 W to achieve 100 Gbps. Thus, the power needed for the RaptorQ code is about 5% of the power needed to read the data off a drive. If other necessary aspects such as network transfer, etc, were incorporated as well, the RaptorQ code share drops well below 5%.

8.7. RaptorQ decode reliability

Extremely high durability is required for storage systems. MDS codes have the property that an object can be recovered from any k fragments for the object. RaptorQ codes are essentially MDS, but not literally MDS: it is not necessarily the case that an object can be recovered from any k fragments. We present simulation results showing the probability of object loss due to RaptorQ codes not being MDS can be essentially ignored.

Consider a liquid system using a (1200, 1000, 200) RaptorQ code, where the target $w = 50$, i.e., the number of available fragments for any object is always at least $k + w = 1050$. If $1/\lambda = 3$ years, then on average there are 400 node failures per year. We group node failures together into consecutive batches of $w - 4 = 46$, and consider testing for decoding using the set of $k + w$ EFIs for available fragments at the beginning of the batch minus the $w - 4 = 46$ EFIs for fragments lost within the batch, and thus we decode from $k + w - (w - 4) = k + 4 = 1004$ EFIs. Note that if decoding is possible from this set of 1004 EFIs then decoding is possible for all objects at each point during the batch.

As shown in Table I, there were no RaptorQ decoding failures in an experiment where 10^{11} random sets of 1004 EFIs were tested for decodability. Since there are on average 400 node failures per year and the batch size is 46, there will be on average $\frac{400}{46} < 10$ such batches per year, and thus these 10^{11} decoding tests cover 10 billion years of simulation time. Thus there would be no decoding failures for RaptorQ within 10 billion years, which is a thousand times a target MTDL of 10 million years.

9. READ REPAIR RATE REGULATION

The simple Poisson process model of node failures discussed in Section 3.1 is useful to benchmark the reliability of a distributed storage system. However, node failures in real systems are not necessarily faithful to a simple statistical model.

Detailed statistical assumptions of the node failure process do not significantly affect reactive repair decisions for a small code system, i.e., the actions taken by the reactive repair process are largely event driven and immediate. For example, repairs are executed as soon as fragments are lost due to node failures, with no dependency on assumptions on future node failures. The anticipated node failure rate is only germane to the initial provisioning of repair resources (e.g., total repair bandwidth).

This is in sharp contrast to lazy repair for a liquid system, where repairs are postponed to achieve greater repair efficiency. The degree to which repairs can be postponed is clearly dependent on some expectation of future node failures. The tradeoff between repair efficiency and MTTDL rests on assumptions about the long term node failure rate. In the liquid system with fixed read repair rate this dependence appears prominently in the choice of the read repair rate. Lower read repair rate results in greater repair efficiency but higher probability of object loss. This is clearly evident in the (tight) MTTDL estimate for the fixed rate system derived in Appendix A.

For example, consider a liquid system using a (402, 268, 134) large code. If the node failure rate λ is $1/3$ per year and the read repair rate is set so that $\lambda \cdot T = 0.21$ then the achieved MTTDL (estimate) is $3.6 \cdot 10^9$ years. If instead λ is 10% higher then the MTTDL is $1.0 \cdot 10^7$ years, and if λ is 20% higher then the MTTDL is $8.7 \cdot 10^4$ years. Thus, the MTTDL of a fixed read repair rate design is sensitive to error in knowledge of λ . If there is uncertainty in the node failure rate then either some repair efficiency must be sacrificed or a reduction in MTTDL must be tolerated for a fixed read repair rate design.

Fortunately, liquid systems can be fitted with an adjustment algorithm that smoothly and automatically adjusts the read repair rate in response to node failure conditions so that accurate a priori knowledge of node failure rates is not required. Suppose for example that the node failure rate increases at some point in time, possibly due to component quality issues. This will result in a relative increase of the number of missing fragments among objects in the repair queue. Moreover, over time the increase in node failure rate can be detected by a node failure rate estimator. Both of these occurrences can signal to the lazy repair process a need to increase the read repair rate. Properly designed, a regulated read repair rate process can achieve significantly larger MTTDL than a fixed read repair rate process running at the same average read repair rate, at the cost of some variation in read repair rate.

Even in the case where the node failure process is a Poisson process of known rate, regulation of the read repair rate can bring significant benefit. To develop a good regulation method it is helpful to consider in more detail the behavior of a fixed read repair rate process. We assume a liquid system with n nodes and a node failure rate λ . We assume a fixed read repair rate that we characterize using T , the time required to repair all objects in the system. We work in the large number of objects limit where the size of a single object is a negligible fraction of all source data. Objects are repaired in a fixed periodic sequence. We view all objects as residing in a repair queue and we use x to denote the relative position in the queue. Thus, an object at position $x = 1$ is under repair and the object in position $x = 0$ has just been repaired. An object in position x will be repaired after a time $(1 - x) \cdot T$ and was repaired a time $x \cdot T$ in the past.

For an object currently in position x let $f(x) \cdot n$ denote the number of its fragments that have been erased due to node failures.² Assume $0 \leq x \leq y \leq 1$ and $0 \leq s \leq t \leq 1$, and that the read repair rate is characterized by T . The number of erased fragments $f(y) \cdot n$ the object will have when

²Note here that while $f(x) \cdot n$ is integer valued we generally consider $f(x)$ as arbitrarily valued on $[0, 1]$. We will make frequent use of this slight abuse of notation.

9 READ REPAIR RATE REGULATION

it reaches position y is given by the transition probability

$$\Pr[f(y) \cdot n = t \cdot n \mid f(x) \cdot n = s \cdot n] = B\left((1-t) \cdot n, (t-s) \cdot n, e^{-\lambda \cdot T \cdot (y-x)}\right) \quad (1)$$

where $B(n, m, q)$ is the probability mass at m of a binomially distributed random variable with parameters (n, q) . That is,

$$B(n, m, q) = \binom{n}{m} \cdot q^{n-m} \cdot \bar{q}^m$$

where we have introduced the notation $\bar{q} = 1 - q$. Applying this to the particular case with initial condition $f(0) \cdot n = 0$ we have $\Pr[f(x) \cdot n = m] = B(n, m, e^{-\lambda \cdot T \cdot x})$. Using Stirling's approximation $\ln(n!) = n \cdot \ln(n) + O(\ln(n))$ we can write $\ln(\Pr[f(x) \cdot n = f]) = -n \cdot E(f, \lambda \cdot T \cdot x) + O(\ln(n))$ where

$$E(f, \lambda \cdot T \cdot x) = -H(f) + \bar{f} \cdot \lambda \cdot T \cdot x - f \cdot \ln(1 - e^{-\lambda \cdot T \cdot x})$$

where $H(f) = -f \cdot \log_2(f) - \bar{f} \cdot \log_2(\bar{f})$ is the entropy function. As a function of f , the exponent $E(f, \lambda \cdot T \cdot x)$ is minimized at $f = 1 - e^{-\lambda \cdot T \cdot x}$, which is also the expected value of $f(x)$, being the solution to the differential equation

$$\frac{d(1 - f(x))}{dx} = -\lambda \cdot T \cdot (1 - f(x))$$

which governs the expected value. Thus, in the large system limit the fraction of erased fragments as function of repair queue position concentrates around the function $1 - e^{-\lambda \cdot T \cdot x}$. In particular we identify $n \cdot (1 - e^{-\lambda \cdot T})$ as the number of fragments that are typically repaired. We can interpret the quantity $f_{\text{tar}} := 1 - e^{-\lambda \cdot T}$ as a *target* fraction of repair fragments. Note that f_{tar} and $\lambda \cdot T$ are in a one-to-one relationship so, given the node failure rate λ , we can view f_{tar} as a specification of the read repair rate. This perspective on specification of system behavior is a convenient one for consideration of regulated read repair rate systems.

To gain a further understanding of the fixed read repair rate system behavior it is helpful to consider typical behavior under object loss. We assume that object loss occurs when an object has more than $r = \beta \cdot n$ fragments erased and consider the number of erased fragments as a function of queue position for such an object. Thus, we consider the distribution $\Pr[f(x) \cdot n \mid f(1) \cdot n > \beta \cdot n]$, which, using Bayes rule, can be written as

$$\begin{aligned} & \Pr[f(x) \cdot n = f \cdot n \mid f(1) \cdot n > \beta \cdot n] \\ = & \frac{\Pr[f(1) \cdot n > \beta \cdot n \mid f(x) \cdot n = f \cdot n] \cdot \Pr[f(x) \cdot n = f \cdot n]}{\Pr[f(1) \cdot n > \beta \cdot n]} \end{aligned}$$

Using (1) and Stirling's approximation as above, we find that in the large n limit the solution concentrates around $f(1) \cdot n = \beta \cdot n$, and, more generally, $f(x) = \frac{\beta}{f_{\text{tar}}} \cdot (1 - e^{-\lambda \cdot T \cdot x})$ which is simply a scaled version of the nominal queue function $1 - e^{-\lambda \cdot T \cdot x}$. Note that this solution satisfies the equation $\frac{d(1-f(x))}{dx} = \lambda \cdot \frac{\beta}{f_{\text{tar}}} \cdot (1 - \frac{f_{\text{tar}}}{\beta} \cdot f(x))$. Thus, for small x , where $f(x)$ is small, the solution corresponds to a typical behavior for the system with node failure rate $\lambda \cdot \frac{\beta}{f_{\text{tar}}}$. Another likely contribution to typical object loss is a skewing of node failures to those nodes which, for the given object in the given repair cycle, have not previously failed. This effect can explain the second factor $1 - \frac{f_{\text{tar}}}{\beta} \cdot f(x)$ which replaces $1 - f(x)$ in the equation for the expected value. Thus, we observe that object loss typically involves a sustained increase in the observed node failure rate over an entire repair cycle and possible skewing of the node failure selection. The fluctuation in node failure rate can be detected by a node failure rate estimator. Both the increased node failure rate and the node selection skewing can be detected from the atypical behavior of the repair queue. We shall

first consider a regulator that assumes a known fixed node failure rate λ and responds only to the repair queue, and later extend the design to incorporate estimators of node failure arrival rates.

Under regulated repair we would expect that if the repair queue state is typical, i.e., the fraction of erased fragments for an object in position x is near its nominal value of $1 - e^{-\lambda \cdot T \cdot x}$ then the read repair rate should be kept near its nominal value. If, however, at some points in the repair queue the number of missing fragments is larger than the nominal amount then the read repair rate should be increased. As a general philosophy we adopt the view that each object in the repair queue takes responsibility for itself to ensure that the probability of object loss for that object is kept small. Each object, knowing its position in the queue and knowing its number of erased fragments, desires a certain system read repair rate such that, if adopted, would adequately protect that object against object loss. The repair system will apply the highest of the desired read repair rates thus satisfying all read repair rate requests for all the objects.

The regulator design defines a function $\phi(f, x)$ such that an object in position x with a fraction f missing fragments requests a read repair rate corresponding to a system repair time of $\lambda^{-1} \phi(f, x)$. Recall that a read repair rate corresponds to time required to repair all object in the system T and that $\lambda \cdot T$ is the expected number of times a single node fails during the time it takes to repair all objects. In this form, $\phi(f, x)$ expresses a desired read repair rate in terms of the desired value for $\lambda \cdot T$. In general $\phi(f, x)$ will be monotonically decreasing in f and increasing in x . Under these assumptions, as well as some assumptions on node failure rate estimators, it follows that the read repair rate will always be determined by certain critical objects in the repair queue. The critical objects are those that were at the tail of the queue when a node failure occurred (see Appendix A for more detail). Note that the fixed read repair rate case corresponds to setting ϕ constant.

As a basis for designing the function $\phi(f, x)$ we adopt the notion that objects desire to recover certain properties of the nominal trajectory. One quantity of particular interest is the probability of object loss for an object. Given a position x in the queue and a fraction $f(x) < \beta$ erased fragments the probability that the object will experience object loss can be expressed using (1). For an object on the nominal trajectory, $f(x) = 1 - e^{-\lambda \cdot T \cdot x} = 1 - (1 - f_{\text{tar}})^x$, this probability is easily seen to be decreasing in x . Thus, if the object enjoys typical behavior early in the repair queue then its probability of object loss decreases. As a design principle, therefore, we can stipulate that an object in position x with a fraction f of erased fragments will request a system repair parameter $\lambda \cdot T$, as represented by the function $\phi(f, x)$, such that the probability of object loss for that object would be the same as an object at that position on the nominal trajectory. Under this stipulation, we obtain $\phi(f, x)$ implicitly as the solution to

$$\sum_{s \cdot n > \beta \cdot n} B(\bar{f}_{\text{tar}}^x \cdot n, (\bar{f}_{\text{tar}}^x - \bar{s}) \cdot n, \bar{f}_{\text{tar}}^{1-x}) = \sum_{s \cdot n > \beta \cdot n} B(\bar{f} \cdot n, (\bar{f} - \bar{s}) \cdot n, e^{-\phi(f, x) \cdot (1-x)}).$$

Note that in this context the quantity $\beta \cdot n$ represents a particular threshold that the regulator attempts avoid reaching and need not necessarily correspond with the erasure correcting limit of the used erasure code. For that reason we introduce a different notation f_T for the target hard threshold. In practice one will typically have $f_T \simeq \beta$.

Simpler expressions of the above principle for deriving appropriate ϕ functions can be obtained in a number of ways. One possibility is note that the above sums are dominated by the term with $s \cdot n = \beta \cdot n$ and then, using the Stirling approximation, to equate the rate functions on both sides. Another approach is to use a Gaussian approximation of the above binomial distributions and equate the integrals that correspond to the above sums. Under the Gaussian approximation we obtain the equation

$$\frac{1 - e^{-(1-x) \cdot \phi(f, x)}}{1 - e^{-(1-x) \cdot \phi_{\text{nom}}}} = \left(\frac{\bar{f}_e - \bar{f}_T}{\bar{f}_{\text{tar}} - \bar{f}_T} \right)^2 \cdot \frac{\bar{f}_{\text{tar}}}{\bar{f}_e} \quad (2)$$

where $\phi_{\text{nom}} = -\ln(1 - f_{\text{tar}})$ and f_e denotes the expected fraction of erased fragments that will be repaired for an object in position x with f erased fragments assuming that the read repair rate

is subsequently fixed according to $\phi(f, x)$, i.e., we have $\bar{f}_e = \bar{f}e^{-\phi(f, x) \cdot (1-x)}$. Substituting this into the above we obtain a quadratic equation for f_e whose solution then determines ϕ . A good approximation is

$$\frac{1 - e^{-(1-x) \cdot \phi(f, x)}}{1 - e^{-(1-x) \cdot \phi_{\text{nom}}}} \simeq \frac{\phi(f, x)}{\phi_{\text{nom}}}$$

under which the level curves of $\phi(f, x)$ correspond exactly to asymptotic trajectories for the regulated system that can be parameterized by their implied target fragment repair fraction f_e . This implies that an object that is controlling the repair rate expects on average to hold the repair constant until it is repaired.

For the simulation results given below we used the Gaussian approximation approach to obtain ϕ . It is not hard to show that ϕ is bounded above (read repair rate never goes to 0). A lower bound is 0 (infinite read repair rate) which is requested if the $f(x)$ is sufficiently close to f_T . In practice the read repair rate will obviously be limited by available resources, and such a limit can be brought into the definition of ϕ by limiting its minimum value and for simulations we typically saturate ϕ at some positive value, often specified as fraction of ϕ_{nom} .

Recall that $\lambda^{-1}\phi$ represents desired time to repair all objects precisely once. This indicates how to incorporate estimates for the node failure rate into the operation of the regulator: one simply replaces λ in the expression $\lambda^{-1}\phi$ with a suitable estimate of λ . In general, an estimator of the node failure rate will be based on past node failures. One pertinent quantity in this context is the appropriate time scale to use to form the estimate. We may choose to vary the time scale as a function of the number of erased fragments the object experiences, or other aspects of the node failure process.

Let us assume a system with a total of \mathcal{O} objects of equal size. Then each object has a position in the queue x with $x \cdot \mathcal{O} \in \{0, \dots, \mathcal{O} - 1\}$. When an object is repaired the position x of each object is incremented by $1/\mathcal{O}$. According to the regulated read repair rate process the repair time for the next object is then given by

$$\frac{\min_x \phi(f(x), x)}{\lambda \cdot \mathcal{O}}$$

where $f(x)$ denotes the fraction of erased fragments for the object in position x . Here we tacitly assume $f(x) < f_T$ for all x but for analytical purposes we may extend the definition of ϕ to all f and x in $[0, 1]$. In practice ϕ is effectively bounded from below by the repair speed limitations of the system. In our simulations we typically assume a floor for ϕ of $\gamma \cdot \phi_{\text{nom}}$ with γ set to a small fraction, e.g. $\gamma = \frac{1}{3}$. Note further that we always have $\phi(0, 0) = \phi_{\text{nom}} = -\ln(1 - f_{\text{tar}})$ and so this gives a lower bound on the read repair rate and we can impose a ceiling on the function ϕ of value ϕ_{nom} without altering system behavior.

The choice of f_{tar} involves a tradeoff: lower values of f_{tar} lead to higher read repair rates, but also a more reliable system, and one that is potentially less prone to large read repair rate variations, whereas a larger value of f_{tar} gives higher repair efficiency at the price of larger read repair rate fluctuations and a reduced object loss safety margin of the system. For simplicity, all the results that we present in this section with the repair adjustment algorithm were obtained by running the algorithm with $f_{\text{tar}} = \frac{2 \cdot x}{3}$, which is a reasonable choice in all but some contrived circumstances.

For certain forms of failure rate estimators it is feasible to compute lower bounds on MTTDL. The details of the models and the calculations can be found in Appendix A. The bounds can be used to designing system parameters such as f_T and more generally ϕ so as to maximize repair efficiency while ensuring a desired MTTDL.

9.1. Regulated read repair rate simulation

Fig. 17 illustrates how the regulated read repair rate automatically varies over time when the node failure rate is $1/3$ years for a liquid system using a (402, 268, 134) large code with storage overhead 33.3%, and the first row of Table II shows a comparison between the regulated repair rate with

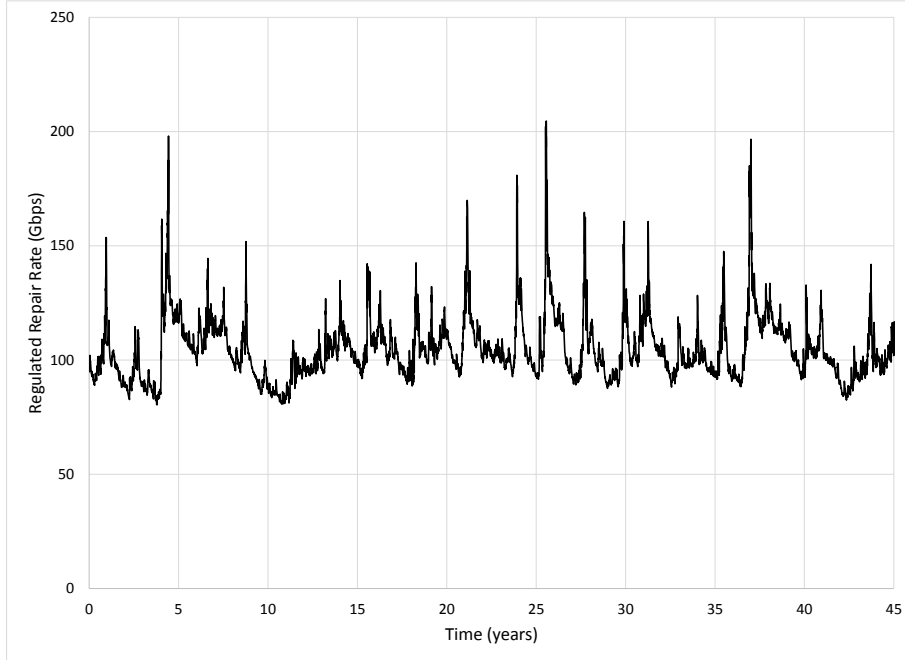


Fig. 17: Repair regulator for a (402, 268, 134) liquid system with fixed $1/\lambda = 3$ years.

an average of $\mathcal{R}_{\text{avg}} = 107$ Gbps and a fixed read repair rate of $\mathcal{R}_{\text{peak}} = 104$ Gbps (set so that the MTDDL is $\approx 10^7$ years). In this table, the MTDDL of $\gg 10^9$ for the regulated read repair rate indicates there were no object losses in 10^9 years of simulation (based on estimates using the results from Appendix A, the MTDDL is orders of magnitude larger than 10^9 years), and $\mathcal{R}_{\text{peak}}$ is a hard upper limit on the used read repair rate. The variability of the regulated read repair rate is moderate, e.g., 99.99% of the time the read repair rate is less than 231 Gbps. The second row of Table II shows a comparison using a (402, 335, 67) large code with storage overhead 16.7%.

System parameters				Regulated read repair rate				Fixed read repair rate		
n	r	β	S	\mathcal{R}_{avg}	$\mathcal{R}_{99\%}$	$\mathcal{R}_{99.99\%}$	$\mathcal{R}_{\text{peak}}$	MTDDL	$\mathcal{R}_{\text{peak}}$	MTDDL
402	134	33.3%	1 PB	106 Gbps	154 Gbps	226 Gbps	311 Gbps	$\gg 10^9$ years	104 Gbps	$3.2 \cdot 10^7$ years
402	67	16.7%	1 PB	298 Gbps	513 Gbps	975 Gbps	1183 Gbps	$\gg 10^9$ years	394 Gbps	$9.0 \cdot 10^7$ years

Table II: Comparison of regulated read repair rate versus fixed read repair rate.

Fig. 18 shows the behavior of regulated read repair rate when the node failure rate is varied as described in Section 5.3. The upper portion of Fig. 18 displays the average node failure rate varying over time, and the lower portion of Fig. 18 displays the corresponding read repair rate chosen by the regulator. This is a stress test of the regulator, which appropriately adjusts the read repair rate to the quickly changing node failure rate to avoid object losses while at the same time using an appropriate amount of bandwidth for the current node failure rate. There were no object losses in 10^9 years of simulation using the regulator running in the conditions depicted in Fig. 18, and based on estimates using the results from Appendix A the MTDDL is orders of magnitude larger than 10^9 years. The results from Appendix A show that the regulator is effective at achieving a large MTDDL when the node failure rate varies even more drastically.

A repair adjustment algorithm could additionally incorporate more advanced aspects. For example, since the bandwidth consumed by the read repair rate is a shared resource, the read repair rate

10 CONCLUSIONS

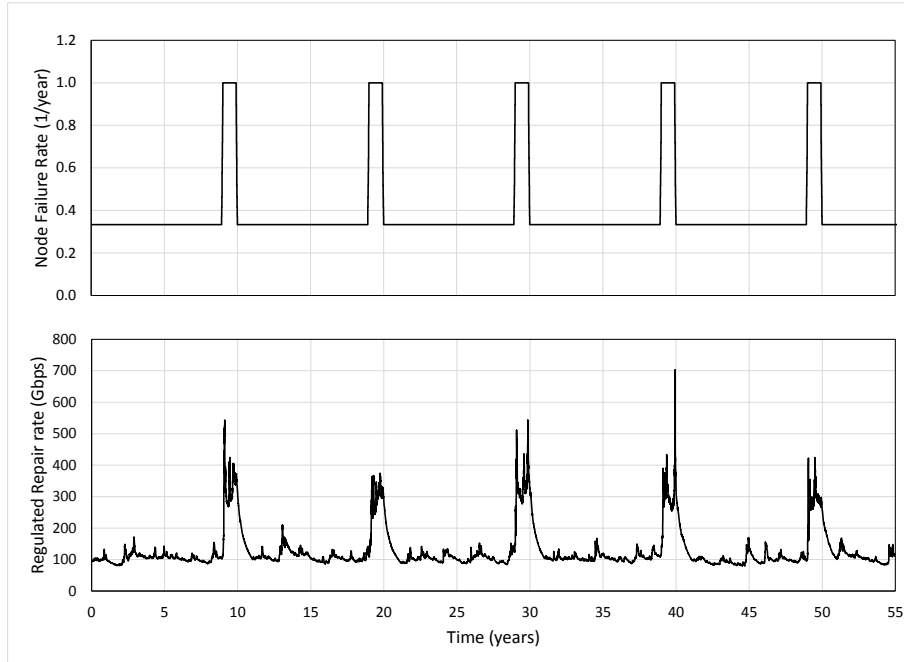


Fig. 18: Regulated repair for a (402, 268, 134) liquid system with node failure rate varying by a factor of three

may also be adjusted according to other bandwidth usage, e.g., the read repair rate can automatically adjust to higher values when object access or storage rates are low, and conversely adjust to lower values when object access or storage rates are high. Furthermore, the read repair rate could be regulated to take into account known future events, such as scheduled upgrades or node replacements.

10. CONCLUSIONS

We introduce a new comprehensive approach to distributed storage, liquid systems, which enable flexible and essentially optimal combinations of storage reliability, storage overhead, repair bandwidth usage, and access performance. The key ingredients of a liquid system are a low complexity large code, a flow storage organization and a lazy repair strategy. The repair regulator design we present provides further robustness to liquid systems against varying and/or unexpected node failure. Our repair and access simulations establish that a liquid system significantly exceeds performance of small code systems in all dimensions, and allows superior operating and performance trade-offs based on specific storage deployment requirements.

We address the practical aspects to take into consideration while implementing a liquid system and provide an example architecture. A liquid system eliminates network and compute hot-spots and the need for urgent repair of failed infrastructure.

While not detailed in the paper, we believe a liquid system provides superior geo-distribution flexibility and is applicable to all kinds of Object, File, Scale-out, Hyper-converged, and HDD/SSD architectures/use-cases. There are also optimizations that the team has worked on that involve pre-generation of EFIs and result in better storage efficiencies at the storage nodes.

Acknowledgment

Rick Dynarski, Jim Kelleman, and Todd Mizenko developed the prototype software for the liquid system described in Section 7, which provides a basic validation of the overall liquid system

properties presented herein. This prototype software is used to cross validate the simulation results presented in Section 5.

Christian Foisy helped to produce the RaptorQ performance results presented in Section 8.5, and also helped cross verify the prototype software for the liquid system. Menucher Menuchehry and Valerie Wentworth provided high level comments that greatly improved the presentation.

We would like to thank Rick, Jim, Todd, Christian, Menucher, and Valerie for their valuable contributions.

A. MTTDL ANALYSIS OF LIQUID SYSTEM

In this section we analyze regulated repair of the liquid system and derive some bounds on MTTDL. In the case of constant repair rate the bound is particularly tight. Although we typically assume that the node failure process is a time invariant Poisson process the analysis can be applied to cases where the rate is not known and time varying. We assume throughout that the number of nodes M is the same as the code length n , and we use n to denote this number.

The framework for the analysis ignores the possibility that an object may be unrepairable. In effect we suppose that as each object leaves the repair queue (because it reaches the head of the queue and, normally, is repaired) a new object enters at the tail. This notion is an analytical convenience that takes advantage of the relative ease of determining or estimating the asymptotic steady state distributions of this system. We assume a large number of objects \mathcal{O} , and for some of the analysis we ignore effects of a finite number of objects and consider the infinite object limit. This has no significant effect on the results but simplifies the analysis considerably.

For the analysis we generally assume that the repair rate is updated upon object repair. In the limit of infinitely many objects this amounts to continuous update. With a finite number of objects the queue state, as sampled at object repair times, the system forms a Markov chain with a unique invariant distribution. In the case where the node failure rate is known it is a finite state chain. Of particular interest is the marginal distribution of the number of erased fragments for the object at the head of the queue. In the case of fixed repair rate we can express this distribution in closed form, and for the regulated case we will derive a bound on the cumulative distribution.

A.1. MTTDL Estimates

We define the MTTDL as the expected time given a perfect initial queue state (all objects have zero erased fragments) until an object achieves $r + 1$ missing fragments prior to completing repair. For the system with node failure rate estimation, we also need to specify the initialization of the estimator. The system is typically designed to achieve a very large MTTDL. We will show in the constant repair rate case that we can obtain a sharp closed-form lower bound on the MTTDL. The bound can be applied even in the case of a mismatch between an assumed node failure rate and an actual node failure rate. In the regulated repair case the analysis is more difficult, but we obtain an efficiently computable generalization of the lower bound for the fixed node failure arrival rate case. We further extend the analysis to cases where the node failure rate is estimated based on the node failure realization.

Between node failure events the system evolves in a purely deterministic fashion; objects simply advance in the repair queue according to the repair rate policy. Under a time invariant Poisson model the node failure interarrival times are independent exponentially distributed random variables and the selection of the failing node is uniformly random over all of the nodes. Thus, the repair can be viewed as a random process consisting of a sequence of samples of the queue state taken upon node failure. The process is generally Markov where the state space consists of a pair of queue states, that just prior to the node failure and that after, and, when there is failure rate estimation, a state representation of the failure rate estimator. The transition to the next state has two components. First, there is the time until the next node failure, which is an exponential random variable. Given that time the queue state advances deterministically from the state after the current node failure to the state just prior to the next. The second component is the selection of the failing node. That choice determines which objects in the queue experience additional fragment loss and which do not. The

update of the state of the failure rate estimators can also occur at this point. Because of the nested nature of the process there will be a point x in the queue such that objects in positions less than x experience fragment loss and those more than x do not, because they were already missing the fragment associated to the failing node. Note that the queue remains monotonic: the set of erased fragments for an object in position x is a subset of the erased fragments for an object in position y whenever $x \leq y$. In the case of node failure rate estimation, those estimates will be updated upon node failure.

A.1.1. The Constant Repair Rate Case. When the repair rate is held constant the deterministic update portion of the process is particularly simple. Let q denote the distribution of the number of erased fragments for the object at the head of the repair queue immediately prior to node failure. Because node failure is a time invariant Poisson process the distribution q is also the distribution of the number of erased fragments for the object at the head of the repair queue at an arbitrary time. In other words, q is the the distribution of the number of erased fragments upon repair under continuous (steady state) operation. We therefore have

$$q(s) = B(n, s, e^{-\lambda T}).$$

Let $F(t)$ denote the number of erased fragments for the object at the head of the repair queue at time t . If t is a node failure time let $F(t-)$ and $F(t+)$ denote value of $F(t)$ immediately prior to and after the node failure. The probability of simultaneous node failures is 0 so we assume $F(t-) \leq F(t+) \leq F(t-) + 1$. Let t be a node failure time, then we have

$$\Pr [F(t+) = r + 1 \wedge F(t-) = r] = \frac{n-r}{n} q(r).$$

Note that this event $F(t+) = r + 1 \wedge F(t-) = r$ corresponds to a data loss event and that the first time this occurs, starting from the perfect queue state, is the first data loss event. We will refer to this event as a *data loss transition*. Thus, the MTTDL is by definition the expected time to the first data loss transition starting from a perfect queue state. Applying Kac's theorem to the Markov chain we have $(\frac{n-r}{n} \cdot q(r))^{-1}$ is the expected number of node failure events between data loss transitions under steady state operation. The expected *time* between such events is given by $(\lambda \cdot (n-r) \cdot q(r))^{-1}$ since $\lambda \cdot n$ is the rate of node failure. To obtain a rigorous bound on MTTDL some further analysis is needed because of differing assumptions on the initial condition. The rigorous correction is negligible in regimes of interest so we use the estimate

$$\frac{1}{\lambda \cdot (n-r) \cdot q(r)} \lesssim \text{MTTDL} \quad (3)$$

as a basis for choosing repair rates (T) as a function of MTTDL.

Starting at time 0 let T_S denote the expected value of the first data loss transition time where the initial queue state is distributed according to the the steady state distribution conditioned on $F(0) \leq r$. Let us also introduce the notation $q(> r) = \sum_{s>r} q(s)$.

LEMMA A.1. *Assume a constant repair rate with system repair time T . Then we have the following bound*

$$\frac{1}{\lambda \cdot (n-r) \cdot q(r)} - \frac{T}{1 - q(> r)} \leq \text{MTTDL}$$

PROOF. We prove the Lemma by establishing the following two inequalities

$$T_S \leq \text{MTTDL} \leq \frac{1}{1 - q(> r)} T + T_S$$

$$\frac{1}{\lambda \cdot (n-r) \cdot q(r)} \leq \frac{1}{1 - q(> r)} T + T_S.$$

Consider a fixed node failure process realization. Among all initial queue conditions with $F(0) \leq r$ the perfect queue will have a maximal time to data loss. Thus, $T_S \leq \text{MTTDL}$.

Assume the system is initialized at time 0 with an arbitrary initial queue state and let us consider the expected time until the next data loss transition. Note that $k \cdot T$ is time when the object initially at the tail of the queue is repaired for the k -th time. Let k^* denote the smallest $k > 0$ where upon repair that object has at most r erased fragments. Then the queue state at time $t = k^* \cdot T$ is the steady state distribution conditioned on the object at the head of the queue having at most r erased fragments, i.e., conditioned on $F(t) \leq r$. Hence, the expected time to the next data loss transition can be upper bounded by $\mathbb{E}(k^* \cdot T) + T_S$. Now, the probability that $k^* = k$ is $(1 - q(> r)) \cdot q(> r)^{k-1}$ hence $\mathbb{E}(k^* \cdot T) \leq \sum_{k>1} (1 - q(> r)) \cdot q(> r)^{k-1} \cdot k \cdot T = \frac{T}{1 - q(> r)}$. It follows that $\text{MTTDL} \leq \frac{T}{1 - q(> r)} + T_S$ and that the expected time between data loss transitions, which is $(\lambda \cdot (n - r) \cdot q(r))^{-1}$, is also upper bounded by $\frac{T}{1 - q(> r)} + T_S$. \square

A.2. Regulated Repair Rate

For the regulated repair rate case the memory introduced into the system by the selection of the maximum requested repair rate makes it quite difficult to analyze. In particular the rate of advancement of the repair queue depends on the entire state of the queue, so there is no convenient way to reduce the size of the Markov chain state space, as was effectively done in the fixed repair rate case. The addition of node failure rate estimation further complicates the problem. It turns out, however, that by approximating queue advancement with a greedy form of operation we can obtain a related lower bound on the MTTDL. When applied to the constant repair rate case the bound is smaller than the one above by a factor of $(1 - f_{\text{arr}})$.

Although we shall consider the case of known node failure arrival rate, much of the analysis is common with the case including node failure estimation. Thus, we will first develop some general results that admit node failure rate estimation.

We introduce the notation $\mathcal{F}_{\text{dis}}(s, t)$ to denote the number of *distinct* node failures in the interval (s, t) . As before, if s is a node failure time we use $s-$ and $s+$ to denote times infinitesimally before or after the node failure respectively. For convenience we assume that $\phi(f, x)$ is bounded below by a positive constant, is defined for all f and x and is non-increasing in f and non-decreasing in x . For computation we may also assume that $\phi(f, x)$ is bounded above by the nominal value $\lambda \cdot T$ so that, in the known node failure rate case, T is an upper bound on the time between repairs for a given object. Recall that since $\lambda \cdot T = \phi(0, 0)$ we always have $\inf_{x \in [0, 1]} \phi(f(x), x) \leq \lambda \cdot T$ so this restriction does not affect the operation of the system. In the case with node failure rate estimation we will typically assume some minimum positive value for the estimated node failure rate so that there will be some finite time T during which a complete system repair is guaranteed.

When we admit failure rate estimation we suppose that each object has its own failure rate estimate. That estimate $\hat{\lambda}$ figures in the requested repair rate through the implied requested system repair time of $\hat{\lambda}^{-1} \cdot \phi(f, x)$. For an object that entered the queue at time s we write its node failure rate estimate at time t as $\hat{\lambda}(s, t)$. The estimate is allowed to depend on past node failures. It can depend on s only through the partition of node failures into ‘past’ and ‘future’ implied by time s . Consequently, all objects that enter the queue between the same pair of successive node failures use the same value for $\hat{\lambda}$. In other words, as a function of s , the estimate $\hat{\lambda}(s, t)$ is piecewise constant in s with discontinuity points only at node failure times. We remark that this can be relaxed to having $\hat{\lambda}(s, t)$ non-decreasing in s between node failure points. In practice the node failure rate estimate used will generally be a function of relatively recent node failure history which can be locally stored. Note that the estimator for an object may depend on its number of erased fragments.

We call an object in the repair queue a *critical* object if it was at the tail of the queue at a node failure point. Letting $s(x, t)$ denote the queue entry time of the object in position x of the queue at time t , an object is critical if $s(x, t)$ is a node failure time. To be more precise, a critical object is

an object that was the last object completely repaired prior to a node failure. Such an object has an erased fragment while it is still at the tail of the queue.

Let $\chi(t)$ denote the set of critical object queue locations $x \in [0, 1]$ at time t . Let $f(x, t) = \mathcal{F}_{\text{Dis}}(s(x, t), t)/n$ and let $\hat{\lambda}_x(t)$ denote $\hat{\lambda}(s(x, t), t)$. Under the current assumptions the repair regulator need only track the critical objects in the system.

LEMMA A.2. *The repair rate is always determined by a critical object, i.e.,*

$$\inf_{x \in [0, 1]} \hat{\lambda}_x^{-1}(t) \cdot \phi(f(x, t), x) = \min_{x \in \chi(t)} \hat{\lambda}_x^{-1}(t) \cdot \phi(f(x, t), x).$$

PROOF. Of all objects that entered the queue between the same pair of successive node failures, the critical object has the minimal value of x and all such objects have the same value of $\hat{\lambda}$ and f . Since $\phi(f, x)$ is non-decreasing in x the result follows. \square

To obtain a bound on the MTTDL we will adopt a greedy approximation in which we effectively assume that objects control the repair rate during their time in the repair queue. To this end, for any time s set $y(s, s) = 0$ and for $t \geq s$ let us define $y(s, t)$ as the solution to

$$\frac{d}{dt} y(s, t) = \hat{\lambda}(s, t) \cdot (\phi(\mathcal{F}_{\text{Dis}}(s, t), y(s, t)))^{-1}$$

which is well defined as long as $y(s, t) \leq 1$. Further, let $\tau(s)$ be the minimal t such that $y(s, t) = 1$. An interpretation of y is that it is the position in the queue of an object that was at the tail of the queue at time s and is designated to be the object that determines the system repair rate as long as it remains in the queue. Let $x(s, t)$ denote the actual position of the object at time t given that it is at the tail of the queue at time s .

LEMMA A.3. *For any s we have $x(s, t) \geq y(s, t)$.*

PROOF. For each fixed f the function $\phi(f, z)$ is non-increasing in $z \in [0, 1]$. Thus if $0 \leq x \leq y \leq 1$ then $(\phi(f, x))^{-1} - (\phi(f, y))^{-1} \geq 0$. Let us fix s and assume $0 \leq y(s, t), x(s, t) \leq 1$. Then we have

$$\begin{aligned} \frac{d}{dt} (x(s, t) - y(s, t)) &= \left(\inf_{z \in [0, 1]} \hat{\lambda}_z^{-1}(t) \cdot \phi(f(z, t), z) \right)^{-1} - (\hat{\lambda}^{-1}(s, t) \cdot \phi(\mathcal{F}_{\text{Dis}}(s, t), y(s, t)))^{-1} \\ &\geq \hat{\lambda}^{-1}(s, t) \left((\phi(\mathcal{F}_{\text{Dis}}(s, t), x(s, t)))^{-1} - (\phi(\mathcal{F}_{\text{Dis}}(s, t), y(s, t)))^{-1} \right) \end{aligned}$$

Thus, $x(s, t) - y(s, t) < 0$ implies $\frac{d}{dt} (x(s, t) - y(s, t)) \geq 0$ and since $x(s, s) = y(s, s) = 0$ we see that $x(t, s) - y(t, s) < 0$ cannot occur by the mean value theorem. \square

For data loss to occur there must be node failure times $s \leq t$ such that both $\mathcal{F}_{\text{Dis}}(s-, t+) \geq r + 1$ and $x(s-, t) < 1$.³ In this event data loss has occurred to the critical object that was at the tail of the queue immediately prior to time s . By Lemma A.3 this implies that $\mathcal{F}_{\text{Dis}}(s-, \tau(s)) \geq r + 1$. The condition $\mathcal{F}_{\text{Dis}}(s-, \tau(s)) \geq r + 1$ denotes a *greedy data loss* where the object arriving at the tail of the queue at time s greedily controls the repair rate and suffers data loss. We refer to the node failure at time s as a greedy data loss node failure. Note that the greedy data loss actually occurs after time s .

A.2.1. The Regulated Repair Case with Fixed Failure Rate. We now assume that the estimated repair rate $\hat{\lambda}$ is constant and may or may not be equal to the actual node failure rate λ . In this case there is no ambiguity in the initialization of $\hat{\lambda}$ and the MTTDL can be lower bounded by the expected value of s where s is the time of the first greedy data loss node failure.

³Here we may interpret $x(s-, t)$ as the limit of $x(s', t)$ as s' approaches s from below.

Assume s is a node failure time and let q now denote the distribution of $\mathcal{F}_{\text{Dis}}(s-, \tau(s-))$. By Kac's theorem, $1/q(> r)$ is the expected number of node failures between greedy data loss node failures. Hence $1/(\lambda \cdot n \cdot q(> r))$ is the expected time between greedy data loss node failures.

Let T_G denote the expected time until the first greedy data loss node failure assuming an unconditioned future node failure process. Note that the time to a greedy data loss is independent of the initial queue state. It follows that $1/(\lambda \cdot n \cdot q(> r))$ is upper bounded by $T + T_G$. This is because if we condition on 0 being a greedy data loss node failure time, then the node failure process after time T is independent of this condition. Recalling that $T_G \leq \text{MTTDL}$ we obtain

$$\frac{1}{\lambda \cdot n \cdot q(> r)} - T \leq \text{MTTDL}$$

In regimes of interest we will have $q(> r) \simeq q(r+1)$ and assuming a fixed repair rate we have $n \cdot q(r+1) = n \cdot B(n-1, r, e^{-\lambda \cdot T}) = e^{\lambda \cdot T} \cdot (n-r) \cdot B(n, r, e^{-\lambda \cdot T})$. If $\hat{\lambda} = \lambda$ then we see that the resulting MTTDL bound is approximately a factor $e^{-\lambda \cdot T} = 1 - f_{\text{tar}}$ worse than the previous one.

Given ϕ it is not difficult to quantize the system to compute the distribution q . To perform the computation we consider quantizing x with spacing δx .⁴ This is essentially equivalent to assuming a finite number, $\mathcal{O} = \delta^{-1}$, of objects and we describe the quantization from that perspective. Assume that repair rates are updated upon object repair. Under regulated repair the repair system then selects the next object repair time, which is now given by the minimum of $\delta \cdot \phi(f(j\delta), j\delta)$ for $j = 0, \dots, \mathcal{O} - 1$. Assume that the next repair time Δ is given and that we are given that the object in position $k \cdot \delta$ has a random number of erased fragments F that is distributed according to $p_k(F)$. Then the distribution of the number of erased fragments for the object when it reaches position $(k+1) \cdot \delta$ is given by

$$p_{k+1}(F) = \sum_{0 \leq F' \leq F} p_k(F') \cdot B(n - F', F - F', e^{-\lambda \cdot \Delta}).$$

To compute the distribution q we proceed similarly but, in effect, track the greedy controlling object as it moves through the repair queue, assuming that that object is the one determining the repair rate. Letting $q_k(F)$ denote the distribution of the number of missing fragments when the object is (greedily) in position $k\delta$. We have

$$q_{k+1}(F) = \sum_{0 \leq F' \leq F} q_k(F') \cdot B(n - F', F - F', e^{-\delta \cdot \phi(F'/n, k\delta)}).$$

Note that the distribution q is simply $q_{\mathcal{O}}(\cdot)$ assuming we initialize with $q_0(F) = \mathbb{1}_{F=1}$ (the distribution with unit mass at $F = 1$.) Using the above expression we see that the distribution q can be computed with complexity that scales as $n \cdot \mathcal{O}$.

Note that if the above computation we instead use the initial condition $\mathbb{1}_F$ then the resulting distribution q provides an upper bound on the steady state distribution of the number of erased fragments upon repair. More precisely, the resulting $q(> s)$ upper bounds the probability that more than s fragments of an object are erased when that object is repaired. Using this value of q we obtain $\sum_{s=0}^n q(> s)$ as an upper bound on the expected number of repaired fragments, which represents the repair efficiency of the system.

In Fig. 19 we plot the computed cumulative distribution of the number of missing fragments upon repair and the computed upper bound. The five shown bounds are different only in the limit placed on the maximum repair rate which is indicated as a factor above the nominal rate. The "1x" case corresponds to the fixed repair rate policy. Already by allowing a "2x" limit a dramatic improvement is observed. It is observed that the greedy bound tracks the true distribution, as found through simulation for the "3x" limited case, with some shift in the number of missing fragments.

⁴An analog of Lemma A.3 can be obtained for the quantized system with the addition of a small correction term, but the impact on the results is negligible for reasonable size systems.

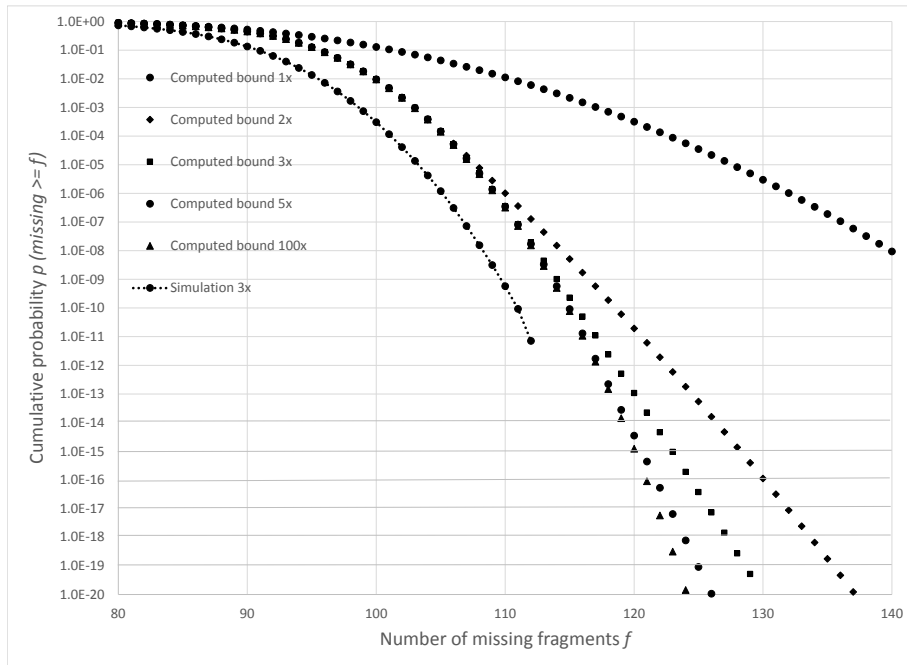


Fig. 19: Comparison of simulated and computed results for regulation with known node failure arrival rate. Different curves place different limits on the peak repair rate, indicated as factors above the nominal rate

Although somewhat conservative, it is evident that the computed bound can provide an effective design tool and that rate regulation can provide dramatic improvements in MTTDL.

Fig. 20 shows how the cumulative distribution of the number of missing fragments changes with node failure rate mismatch. We set the repair rate limit to 10 times the nominal rate in increased the node failure rate by factors of 1.25, 1.5, and 2.0 respectively. The regulator responds by increasing the average repair rate and reaching a new asymptotic queue state equilibrium with respectively larger repair targets. This is indicated in the shifts of the cumulative distribution. For very large shifts in failure rate this level of adaptation would be inadequate to ensure data durability. Fortunately node failure rate estimation is easily incorporated into the regulator design, leading to a significantly more robust system.

A.2.2. The Regulated Repair Case with Failure Rate Estimation. The rate regulator as described above is designed to concentrate the repair queue near its nominal behavior when the actual node failure rate matches the design's assumed failure rate. As demonstrated, the regulator can compensate for some variation in the failure rate. If the actual node failure rate is persistently different from the design assumed rate then the regulator typical behavior is shifted relative to the nominal one. By using an estimator of the node failure rate the system can bring the typical behavior back toward the nominal one. Moreover, by using an estimate of the node failure rate the system can adapt to much larger variations in the actual node failure rate.

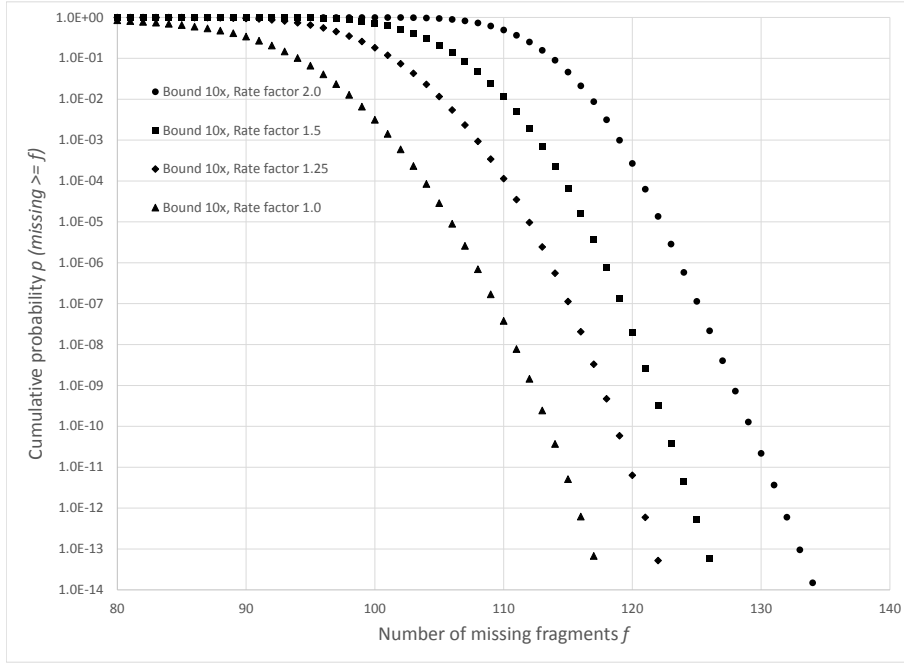


Fig. 20: Cumulative distribution bounds with mismatch between assumed design node failure rate and actual node failure rate.

Many forms of estimators may be considered. We shall consider some relatively simple forms that are amenable to computation analysis similar to the fixed arrival rate case. These estimates are based on first order filtering of some function of the node failure interarrival times. A typical update of the node failure interarrival time will take the form

$$\xi(\hat{T}) \leftarrow \alpha \cdot \xi(\hat{T}) + \bar{\alpha} \cdot \xi(T).$$

where T is a most recent node failure interarrival time and the functional ξ is an invertible function. A convenient case to consider is $\xi(T) = T$ but other potentially interesting forms are $\xi(T) = \log T$ or $\xi(T) = 1/T$, or regularized versions of these functions such as $\xi(T) = \log(1+T)$. For simplicity we will focus on the case $\xi(T) = T$. The node failure rate estimate $\hat{\lambda}$ is simply given as \hat{T}^{-1} .

Corresponding to our assumptions on failure rate estimators, we assume that the value of α used by an object is a function only of the number of erased fragments for that object and the total number of node failures that have occurred since that object entered the repair queue, for which we introduce the notation $\mathcal{F}_{\text{All}}(s, t)$. Let us now consider how to compute the distribution of the greedy behavior. In this case we take a slightly different approach than previously. We aim to recursively compute the joint greedy queue position and failure rate estimate distribution at node failure times. Since a node failure may or may not result in a fragment erasure for the greedy object, we must also jointly consider the distribution of the number of erased fragments.

Let us fix s and consider $F(t) = \mathcal{F}_{\text{Dis}}(s, t)$, $F^+(t) = \mathcal{F}_{\text{All}}(s, t)$, $\hat{\lambda}(s, t)$, $y(s, t)$, for $t \leq \tau(s)$ and extend these definitions for $t \geq \tau(s)$. Specifically, for $t \geq \tau(s)$ set $y(s, t) = 1$, set $F(t) = F(\tau(s))$, set $\hat{\lambda}(s, t) = \hat{\lambda}(s, \tau(s))$. We will leave $F^+(t) = \mathcal{F}_{\text{All}}(s, t)$ unchanged, i.e. $F^+(t)$ continues to count node failures. In short we track $F(t)$, $y(s, t)$, and $\hat{\lambda}(s, t)$ until $y(s, t) = 1$ which then marks an absorbing state.

We will now briefly discuss how, given F^+ , we can computationally update a joint distribution of F , y , $\hat{\lambda}$ to the next node failure point. Assume that at time t_0 the greedy object has an estimate $\hat{\lambda}(s, t_0)$ and greedy queue position $y = y(s, t_0)$ and a certain value of F . Then the joint distribution of y and $\hat{\lambda}(s, t)$ at the next node failure point lies on a curve parameterized by t_1 , the next node failure time. Note that since we assume that no additional node failures occur, we have that $\hat{\lambda}(s, t)$, $F(t)$, and $F^+(t)$ are all constant and we can drop the dependence on s and t . Note also that α is determined and is independent of the next node failure time. For $t \geq t_0$ such that $y(s, t) < 1$ we have

$$y(s, t) = y(s, t_0) + \hat{\lambda} \int_{t_0}^t \phi^{-1}(F/n, u) du.$$

Assuming $t = t_1$ is the next node failure the update of $\hat{\lambda}$ is given by

$$\hat{T}(s, t) = \alpha \cdot \hat{T}(s, t_0) + \bar{\alpha}(t - t_0)$$

and we obtain the curve $y(s, t_1)$, $\hat{T}(s, t_1)$ for $y(s, t_1) < 1$. Note, however, that at $y(s, t_1) = 1$ we actually have the discontinuity $\hat{T}(s, t_1) = \hat{T}(s, t_0)$ since in this case the object reached the head of the queue prior to additional node failures. To obtain the probability density on this curve observe that until the time of the next node failure t_1 we have

$$t - t_0 = \hat{\lambda}^{-1} \int_{y(s, t_0)}^{y(s, t)} \phi(F/n, u) du.$$

The distribution of $t_1 - t_0$ is exponential with rate $\lambda \cdot n$ so the probability that $t - t_0$ is greater than a quantity $t' > 0$ is given by $e^{-\lambda \cdot n \cdot t'}$. Since $y(s, t)$ is monotonic in t , the probability that $y(s, t_1)$ is greater than z where $y(s, t_0) \leq z < 1$ is given by

$$\Pr[y(s, t_1) > z] = e^{-\lambda \cdot n \cdot \hat{\lambda}^{-1} \int_{y(s, t_0)}^z \phi(F/n, u) du}$$

and this gives the distribution along the curve $y(s, t_1)$, $\hat{T}(s, t_1)$.

Let us introduce the notation $Q_{F^+}(F, y, \hat{\lambda})$ to denote the conditional density of the extended definitions of F , y , $\hat{\lambda}$ conditioned on F^+ . Using the above described calculations we can update this distribution to the distribution at the next node failure point. We will use the notation $\tilde{Q}_{F^+}(y, \hat{\lambda})$ to denote this updated distribution. Assuming $y < 1$, the next node failure is incrementing with probability $(1 - F/n)$ thus for $y < 1$ we have

$$Q_{F^+}(F, y, \hat{\lambda}) = (F/n)\tilde{Q}_{F^+-1}(F, y, \hat{\lambda}) + (1 - (F - 1)/n)\tilde{Q}_{F^+-1}(F - 1, y, \hat{\lambda}).$$

and for $y = 1$ we have

$$Q_{F^+}(F, y, \hat{\lambda}) = \tilde{Q}_{F^+-1}(F, y, \hat{\lambda}).$$

With suitable quantization of y and $\hat{\lambda}$ the above procedure can be used to compute the greedy distribution.

To obtain the greedy bound for an arbitrary initial time we initialize the distribution Q_0 so that the marginal distribution of y is $\mathbb{1}_y$ (and, necessarily, $F = 0$ with probability 1.) Given a Poisson node failure model and a particular estimator the distribution $\hat{\lambda}(s, s)$ may be determined. In general

the distribution of $\hat{\lambda}$ can be set according to the desired assumptions of the computation. The limit as F^+ tends to $+\infty$ of $Q_{F^+}(F, 1, \hat{\lambda})$ is the greedy approximation to the distribution upon repair of F and $\hat{\lambda}$. The probability of large F^+ with $y < 1$ decays exponentially in F^+ , so this computation converges quickly in practice.

To obtain a bound on the MTTDL we need to perform the corresponding computation for critical objects. This requires only one additional initialization step in which the distribution $Q_1(F, y, \hat{\lambda})$ is modified to project out y and put all probability mass on $y = 0$. Note that when $F^+ = 1$ we have $F = 1$ with probability 1. With such a modification of the initialization process we can obtain a lower bound on the MTTDL as before. We remark, however, that in this case we take the definition of MTTDL to include a random initialization of $\hat{\lambda}$ according to the distribution of $\lambda(s, s)$ and the above computation should also use this initial condition.

Fig. 21 shows the cumulative distribution of the number of missing fragments upon repair q as obtained from simulation and the computed bound. The simulation objects used for an estimate of node failure rate the average of the node failure interarrival times for the last $155 - F$ node failures, where F denotes the number of erased fragments in a object. Note that the value of 155 was chosen slightly larger than $r + 1 = 135$. The computed bound used a slightly different estimator consistent with the above analysis. The estimator used is a first order filter but, to best model the corresponding simulation, the time constant of the filter was adjusted as a function of F to match $155 - F$. In both cases the peak repair rate was limited to 3 times the nominal rate as determined by the actual node failure rate. It can be observed from these the plots that even compared with know node failure arrival rate, the estimated case can give larger MTTDL. This is intuitive since data loss is generally associated to increases in node failure rate and the increase in the estimated node failure rate case gives rise to faster repair. In addition, as shown previously, a regulator with rate estimation can accomodate large swings in node failure rates.

REFERENCES

- [1] M. Belshe, R. Peon, M. Thomson. IETF RFC7540. Hypertext Transfer Protocol Version 2 (HTTP/2). Internet Engineering Task Force, May 2015.
- [2] J. Bloemer, M. Kalfane, M. Karpinski, R. Karp, M. Luby, D. Zuckerman. An XOR-Based Erasure-Resilient Coding Scheme. *ICSI Technical Report*, No. TR-95-048, August 1995.
- [3] B. Calder, J. Wang, A. Ogun, N. Nilakantan, A. Skjolsvold, S. McKelvie, Y. Xu, S. Srivastav, J. Wu, H. Simitci, J. Haridas, C. Uddaraju, H. Khatri, A. Edwards, V. Bedekar, S. Mainali, R. Abbasi, A. Agarwal, M. Fahim ul Haq, M. Ikram ul Haq, D. Bhardwaj, S. Dayanand, A. Adusumilli, M. McNett, S. Sankaran, K. Manivannan, L. Rigas. Windows Azure Storage: A Highly Available Cloud Storage Service with Strong Consistency. *SOSP '11*. Oct. 23-26, 2011, Cascais, Portugal.
- [4] J. Cowling. <https://code.facebook.com/posts/253562281667886/data-scale-june-2016-recap/>
- [5] Y. Chen, R. Griffith, D. Zats, A. D. Joseph, R. Katz. Understanding TCP Incast and Its Implications for Big Data Workloads *University of California at Berkeley, Technical Report, 2012*
- [6] A. Dimakis, P. Godfrey, Y. Wu, M. Wainwright, K. Ramchandran. Network coding for distributed storage systems. *IEEE Infocom*, May 2007
- [7] A. Dimakis, P. Godfrey, Y. Wu, M. Wainwright, K. Ramchandran. Network coding for distributed storage systems. *IEEE Transactions on Information Theory*, Vol. 56, No. 9, pp. 4539-4551, September 2010
- [8] D. Ford, F. Labelle, F. Popovici, M. Stokely, V. Truong, L. Barroso, C. Grimes, and S. Quinlan. Availability in globally distributed storage systems. *Proceedings of the 9th USENIX Symposium on Operating Systems Designs and Implementation*, pp. 1-7 Oct. 2010.
- [9] P. Gopalan, C. Huang, H. Simitci, S. Yekhanin. On the Locality of Codeword Symbols. *IEEE Trans. Inf. Theory*, vol. 58, No. 11, pp. 6925-6934, Nov. 2012.
- [10] C. Huang, H. Simitci, Y. Xu, A. Ogun, B. Calder, P. Gopalan, J. Li, and S. Yekhanin. Erasure Coding in Windows Azure System. *USENIX Annual Tech. Conference*, Boston, MA, 2012.
- [11] G. Joshi, Y. Liu, E. Soljanin. Coding for fast content download. *Communication, Control, and Computing (Allerton)*, 50th Annual Allerton Conference, IEEE, pp. 326-333, Oct. 2012.
- [12] J. Lacan, V. Roca, J. Peltotalo, S. Peltotalo. IETF RFC5510. Reed-Solomon Forward Error Correction (FEC) Schemes. Internet Engineering Task Force, April 2009.

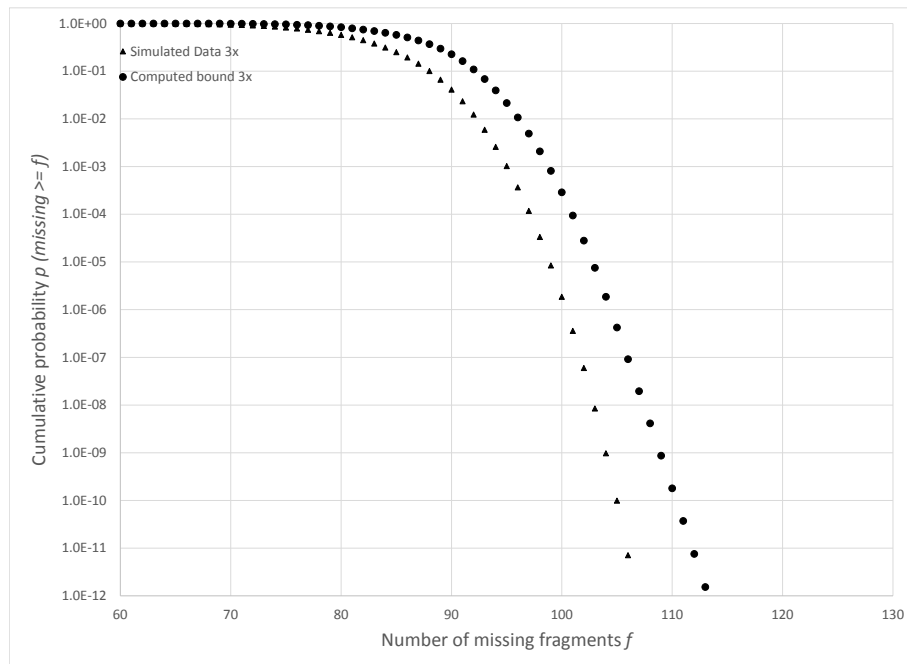


Fig. 21: Comparison of simulated and computed bound for repair regulation with failure rate estimation and a peak repair rate limit of 3 times the nominal rate.

- [13] M. Luby, A. Shokrollahi, M. Watson, T. Stockhammer, L. Minder. IETF RFC6330. RaptorQ Forward Error Correction Scheme for Object Delivery. Internet Engineering Task Force, August 2011.
- [14] M. Luby. Capacity Bounds for Distribute Storage. *submitted to IEEE Transactions on Information Theory*, October 2016.
- [15] S. Muralidhar, W. Lloyd, S. Roy, C. Hill, E. Lin, W. Liu, S. Pan, S. Shankar, V. Sivakumar, L. Tang, S. Kumar. f4: Facebook's warm BLOB storage system. *11th USENIX conference on Operating Systems Design and Implementation*, pp. 383-398, 2014.
- [16] R. Recio, B. Metzler, P. Culley, J. Hilland, D. Garcia IETF RFC5040 A Remote Direct Memory Access Protocol Specification Internet Engineering Task Force, October 2007.
- [17] L. Rizzo. Effective erasure codes for reliable computer communication protocols. *ACM SIGCOMM computer communication review*, Vol. 27, No. 2, pp. 24-36, April 1997.
- [18] R. Rodrigues, B. Liskov. High Availability in DHTs: Erasure Coding vs. Replication. *Peer-to-Peer Systems IV*, pp. 226-239, Springer Berlin Heidelberg, 2005
- [19] M. Sathiamoorthy, M. Asteris, D. Papailiopoulos, A. Dimakis, R. Vadali, S. Chen, and D. Borthakur. XORing Elephants: Novel Erasure Codes for Big Data. *Proceedings of the VLDB Endowment*, pp. 325-336 Vol. 6, No. 5, 2013.
- [20] A. Shokrollahi, M. Luby. Raptor Codes. *Foundations and Trends in Communications and Information Theory*, 2011, Vol. 6: No 3-4, pp 213-322.
- [21] H. Weatherspoon, Hakim, J. Kubiatowicz. Erasure coding vs. replication: A quantitative comparison. *Peer-to-Peer Systems*, pp. 328-337, Springer-Verlag, 2002.
- [22] <http://docs.ceph.com/docs/master/rados/configuration/pool-pg-config-ref/>
- [23] <https://github.com/google/snappy/blob/master/README>
- [24] SM863a Specification Sheet <http://www.samsung.com/semiconductor/minisite/ssd/product/enterprise/sm863a.html>

REFERENCES

- [25] Hard Drive Data and Stats <https://www.backblaze.com/b2/hard-drive-test-data.html>

REFERENCES