

Providing Both Scale and Security through a Single Core Probabilistic Protocol*

Ramsés Morales and Indranil Gupta
University of Illinois at Urbana-Champaign
Department of Computer Science
Urbana, IL 61801
{rvmorale,indy}@cs.uiuc.edu

Abstract

Distributed systems are typically designed for scale and performance first, which makes it difficult to add security later without affecting the original properties. This paper proposes the design of the Folklore persistent distributed storage system, which adopts an alternative design methodology. Folklore’s design relies on a single core protocol for providing both probabilistic scalability and untraceability, the latter being a special notion of probabilistic security. The core protocol is a biologically inspired model of endemic replication that migrates replicas of files among all hosts in a continuous and proactive manner. The emergent behavior is chaotic, meaning that the exact number and location of all replicas of any file is changing all the time. This makes it difficult for an attacker to target any file. Yet, the protocol is scalable – it consumes constant per-host bandwidth, and the number of replicas per file stays close to a small self-stabilizing value. The self-stabilizing value is reached even if only one replica survives a massive attack. The simplicity of the core protocol allows augmentation with mechanisms that allow data integrity, availability, and updatability. We describe the internals of the Folklore system, present attack analysis, and give experimental results from a prototype that shows high resilience to large-scale attacks.

1. Introduction

Several networked information stores such as digital libraries [10], global-scale storage [9], “brick”-based storage [13], etc., have emerged over the past few years. All require an ability to scale with thousands of hosts, as well as to provide such notions of security as resistance to attackers, fail-

ures, and churn. However, many existing mechanisms that ensure security tend to be unscalable. For instance, Byzantine agreement protocols, e.g., [5] require several rounds of all-to-all communication, preventing scalability with even hundreds of hosts.

We are exploring the design of a persistent and distributed file store called *Folklore* that provides both scalability and certain notions of security. This is achieved by relying on a *single* probabilistic core protocol. This core, called the *endemic replication protocol* [8], continuously migrates replicas of files among all the hosts in the system. The protocol is inspired by phenomena resulting in survival of diseases, e.g., influenza.

The **basic notion of security** provided by Folklore is called **untraceability**. The continuous migration of replicas due to the endemic replication occurs *chaotically*, i.e., the exact number and location of the replicas of a given file is changing all the time, and in an unpredictable manner. This makes it very difficult for any attacker to discover and destroy all replicas of a targeted file. In fact, the average attack window times are very short (typically tens of seconds), beyond which the endemic protocol automatically creates additional replicas that the attacker will not know about.

Folklore is scalable because in spite of the chaotic nature of the replication, the number of replicas of each file stays centered around a stable value, resulting in a small storage utilization. The bandwidth utilization at each Folklore host is surprisingly low – it is independent of the number of hosts in the system, and depends only on the average file size. For an average file size of 128KB, our experiments indicate that the per-host bandwidth is 3.38 kBps.

The design of Folklore augments the core protocol with simple mechanisms that provide (1) **data availability and updatability**, i.e., the owner of a file can quickly access and modify the file, (2) **data integrity**, i.e., a file cannot be corrupted by an attacker.

*This research was supported in part by NSF CAREER grant CNS-0448246.

LOCKSS [10] is the only other system that provides probabilistic security and scale. Several persistent distributed stores have been built using distributed hash tables (e.g., [9, 12]), but these do not provide untraceability. In those systems, a surreptitious attacker can first find out the IP addresses of the replica nodes for a given file (through normal lookups), and then bring down these computer hosts (either manually or through a network-based attack), thus removing the file from the system. This attack also works when erasure-coding based replication strategies [3] are used.

2. Design of Folklore

2.1. Core Protocol: Endemic Replication

We describe the variant of the endemic replication protocol [8] that is used in the Folklore system. Consider a host (henceforth a “node”) H that is currently storing replicas of some of the files inserted into the system. We say that the node is *stashing* each of these files, or that it is in a *stash* state w.r.t. each of these files. For all the files in the system that are not locally stashed at H , H is either *receptive* or *averse*. H is said to be receptive to a file f when it is willing to accept a transfer of f 's replica from another stasher node of f . If it refuses any such transfer, it is said to be averse w.r.t. that file. H maintains a list of names for files that it is averse to, while names of the remaining files (for which it is receptive) are not remembered.

The actions of the endemic replication protocol are determined by four protocol parameters - protocol period T , fixed probabilities $\alpha, \gamma \in [0, 1]$, and an even integer β . The basic endemic protocol is periodic, with each node executing certain actions once every T time units. Protocol periods are asynchronous across nodes. Parameter values are fixed at all nodes.

At the start of each protocol period at node H , several actions are executed: (1) First, each locally stashed file is moved into the averse state with probability γ . In the averse state, the file can be deleted, although its identifier needs to be remembered in order to refuse any transfers. (2) Then, each averse file is forgotten (i.e., moved into receptive state) with probability α . (3) Finally, H creates a constant sized *stash advertisement* message, containing identifiers of some locally stashed files. The *stash advertisement* is sent to a constant number $\frac{\beta}{2}$ of target nodes, each one selected uniformly at random from across the group¹. Target nodes receptive to the advertised files may request a transfer from H . In turn, H requests these target nodes to advertise their stashed files, and transfers to itself some files that it is receptive to.

¹We use a variant of the SWIM membership protocol [6] at each node to maintain partial, local, membership lists.

The analysis of [8] showed that for any given file, the fractions of nodes that are (*receptive, stasher, averse*) stabilize around a single equilibrium point - $(\frac{\gamma}{\beta}, \frac{1-\gamma}{1+\frac{\alpha}{\beta}}, \frac{1-\gamma}{1+\frac{\alpha}{\beta}})$. This equilibrium point is stable, which means the endemic protocol always self-stabilizes around it. However, the protocol exhibits chaotic perturbation in the vicinity of this point because of random and unpredictable choices by each node while applying the probabilistic protocol actions (1-3) above.

2.2. Basic File Operations

A file inserted into the Folklore file system is assigned a globally unique *id* using the base64-encoded SHA1 [2] hash of the owning node's IP address and the filepath at the owning node.

File replicas are encrypted for privacy and plausible deniability to stashers. We use symmetric encryption that is calculated per 4 KB file-blocks. This partial encryption will help in propagating partial file updates (Section 2.4). We use the efficient Blowfish [1] encryption algorithm in Cypher Feedback (CB) mode, where each cleartext 4KB block will result in a 4KB cipher.

File replicas are provided with a signature, computed as $\text{privateKey.sign(SHA1(file, version, id))}$. A node will refuse to store a file replica whose signature cannot be validated. The signature is verified with an *authentication token*, a third party signature of the file owner's public key concatenated with its IP address. The public key of the third party is preinstalled on all nodes. Such a third party is unavoidable due to the Sybil attack impossibility result [7].

Finally, for a newly inserted file, the owning node warms up the system by forcing a small random set of nodes to become stashers for the new file; other nodes are of course receptive.

2.3. Stash-Receptive Contacts

As explained before, each node i periodically sends out a *stash advertisement* message to target node(s) j , chosen uniformly at random. The message interchange is depicted in Figure 1, and detailed below.

Each file in all advertisement messages is specified as $\langle id, version \rangle$, the latter field being the owner-set version, incremented whenever the owner makes an update. Suppose node i sends a *stash advertisement* message, and the recipient node is j . j replies to i with a *stash-receptive* message, which contains two sets of filenames: a *stash set* and a *receptive set*.

The *stash set* advertises all files stashed at node j that were not specified in the *stash advertisement* message (and could thus be transferred to i). The *receptive set* only identifies those files chosen from the *stash advertisement* mes-

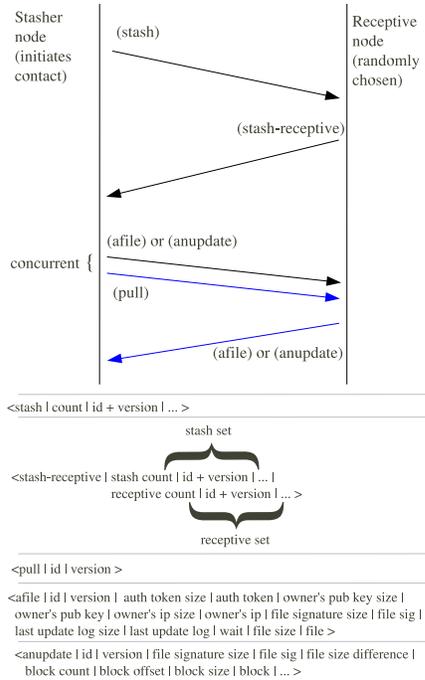


Figure 1. Message interchange when a node contacts a peer with its stash advertisement.

sage from i that either (1) j is receptive to, or (2) that are also stashed at j but have a lower version number.

To save bandwidth, we allow a node to retain a file replica even when it becomes averse or receptive to the file. Accordingly, condition (1) for selecting the receptive set changes to include those files that: (1)* j is receptive to, and if it is retaining a copy of, that has a lower version number than advertised by i . In addition, any locally retained files at j , that are advertised with a lower version number in i 's stash advertisement, can be moved into a stash state at node j .

Push-Pull mechanism: On receiving the *stash receptive* reply M' from node j , node i selects *one* id from M' 's *receptive set* that is still stashed at i . Node i then *pushes* this file id to j . If the difference between version numbers of this file locally and in M' is 1, a per-block update (*anupdate*) is done, otherwise a full file transfer (*afile*) is used.

Node i also selects one of the files in the *stash set* of M' that is still receptive at node i . Node i sends a *pull* message, containing this file's id , to node j , which then replies with either a file transfer or update. For all remaining id 's in the *stash set* of M' that are retained locally at node i , those files with local version numbers that are greater or equal are quietly moved into the stash state; the ones with lower local

version numbers are deleted.

To ensure that an inserted file ramps up to the stable number of stashes, initial replicas are associated with a *TTL=time to live*. Any replica created with non-zero TTL is forced to stay stashed for the next TTL periods, after which normal Folklore actions apply. TTL is initialized by the owner, as described below. Every new replica created by a *stash-receptive* contact inherits the TTL from its stashed parent, less a quantity $w = \frac{S_i}{\frac{\beta}{2}}$. Here, S_i is parent node's stash size, thus w is the expected time (in periods) to advertise the replica once. The owner initializes the TTL to $w \cdot \lceil \log_{(\frac{\beta}{2}+1)} S_\infty \rceil$ periods. S_∞ is an estimate of the stable number of stashes, e.g., based on expected system size. The total number of forced replicas grows by a factor of $(\frac{\beta}{2}+1)$ every w periods, thus the logarithm's base. It should be evident that all TTLs reach 0 a while after file insertion.

Each node expects on average β contacts every protocol period, and each contact will potentially result in two transfers. Therefore, the theoretical bandwidth used at each node is $\frac{2\beta * averageFileSize}{T}$. Our optimizations reduce the bandwidth at each node by 20.84%.²

2.4. File Location, Update, and Recovery

The untrusted nature of hosts allows us to impose the restriction that all such accesses are done by the owner node of the file³. This allows us to build in relatively simple solutions that leverage off a *heartbeating strategy*.

The heartbeating strategy requires each file f to be associated with a heartbeat probability parameter $\epsilon(f) \in [0, 1]$. At the beginning of a protocol period at node H , for each locally stashed file f , and with probability $\epsilon(f)$, H sends its owner a heartbeat message specifying the file's $\langle id, version \rangle$. The expected number of heartbeat messages received at a file's owner, per protocol period, is thus $hb = \epsilon(f) \times N \times \frac{1-\frac{\gamma}{\beta}}{1+\frac{\alpha}{\beta}}$, where N is the number of non-faulty nodes currently in the system. If the expected range for N is known, ϵ can be set so that the owner receives a constant number of heartbeats per file per period. ϵ also trades between response time for a lookup and the bandwidth used by heartbeats.

Due to heartbeats, the owner of a file f is not required to store the file anymore. Only knowledge of the latest version number and symmetric encryption key suffices. When the owner wishes to locate a copy, it simply waits to receive the next heartbeat for file f with the latest version number.

²For comparison, traces used in [12] had a mean file size of 10KB, and [11] found that only 7% of files are bigger than 96KB on an NT file system. New research addresses how to keep the replica size constant, independent of file size.

³A new research project based on the core endemic protocol supports update from multiple users

The heartbeating mechanism can be used to spread updates. An update can be spread immediately to a node that sent a heartbeat – if the version number difference is 1, then only modified blocks are transferred, otherwise the entire new file is transferred.

Suppose the stable number of stashes for a file is S_∞ , and the average number of heartbeats received per protocol period is hb . Then, the number of updated stashes S_i after i rounds since the update start satisfies $(S_\infty - S_i) = (S_\infty - S_{i-1})[(1 - \frac{hb}{S_\infty})(1 - \frac{\beta \cdot S_{i-1}}{2N})(1 - \frac{1}{N})^{\beta \cdot S_{i-1}/2}]$. Analysis and experiments show that within the square brackets, the last two terms (stash-stash and stash-receptive contact) dominate the first term (heartbeating).

3. Attack Analysis of Folklore

We discuss vulnerabilities of the Folklore system and its behavior under various types of attacks. Here –and in the experiments– we assume an attacker with enough resources to compromise as much nodes as he wants at any time.

Replica or Version Corruption: A malicious stasher or a man in the middle that modifies a replica, or its version, or its id will be detected by a good node during signature verification. This can also be used to raise an alarm about malicious nodes.

Owner Node Takeover: If the owner node is taken over, but the private key of the owner is not compromised, the attacker can view heartbeats and files, but will be unable to propagate updates. Thus, a Folklore file cannot be corrupted by an owner takeover attack. To decrease the chance of compromising the private key, it could be stored on secure media and away from the owner node when not needed.

Intentional File id Collisions: If a node detects an *id* conflict, it will ask all conflicting file owners to rename their files. If any of the owners refuse, they are blacklisted at the detecting node.

Sybil Attack: This attack occurs when a large number of nodes are (or appear to be) malicious in the system. For instance, an extreme variety works by breaching a gateway. Although such an attack can cripple any system, Folklore allows a good node to be able to detect such an attack if the number of heartbeats received per period is too low.

4. Experimental Results

Folklore is implemented in C. We present experimental results from a simulated Folklore system running on a 3GHz PIV, 1GB RAM, Linux PC. The system has 1000 nodes, with protocol period $T = 5$ minutes, and $\alpha = 0.05$, $\gamma = 0.4$. A node contacts 5 peers each protocol period (thus $\beta = 10$).

Chaotic Replication: At time $t = 10000$ s, 1000 files are simultaneously inserted into the system, each from a unique

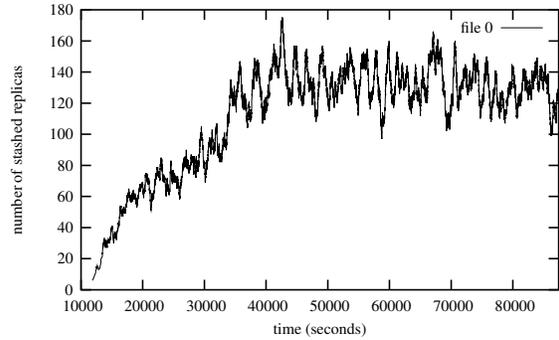


Figure 2. Replica propagation. 1000 files are inserted into the network.

owning node. Figure 2 shows how the number of replicas (for file 0) varies over time; similar trends were observed for the other files. The stable expected number of replicas is 120. The chaotic behavior of the protocol around this stable value is evident from Figure 2.

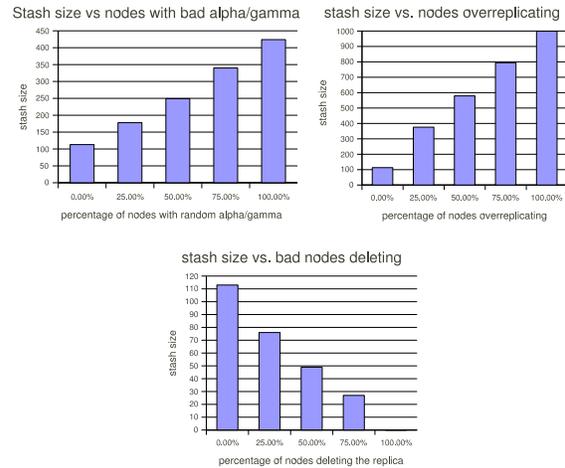


Figure 3. Effect of Parameter Attack, Over-replication and Replica Destruction.

Parameter Attack: Figures 3(a)-3(c) show the effect of a large collection of malicious nodes using bad values for protocol parameters. Figure 3(a) shows the effect of a **random parameter attack** - a fraction of nodes choose random values for $\alpha, \gamma \in [0, 1]$. With 50% of the nodes malicious, the number of replicas of each file increases only two-fold. Figure 3(b) shows the effect of an **over-replicating attack**. Over-replicating nodes reduce the protocol period to 30 s and lock the replica in stash state, i.e., never changes to averse. With 25% of malicious nodes, the number of replicas only increases by a factor of 3. Figure 3(c) shows a

deletion attack. Even when 75% nodes delete all files they purport to stash, the number of replicas stays non-zero.

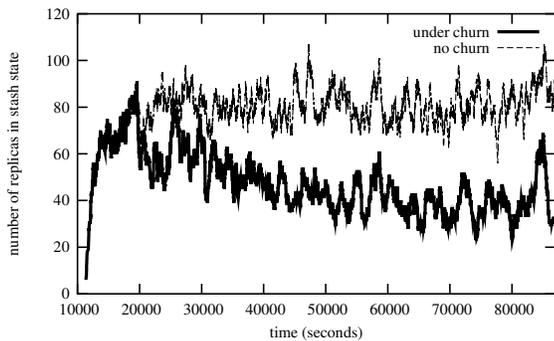


Figure 4. Effect of Churn Attack.

Churn Attack: In a network of 526 nodes and 526 files, churn was introduced at time $t = 18557$, using churn traces from the Overnet network [4]. The rest of the protocol parameters remained the same. Figure 4 shows the typical effect for a file – churn causes the number of replicas to halve, but not disappear. During this time-span, around 70% of the hosts have an availability of at most 0.5 [4]. The protocol finds a new equilibrium point for the number of stash replicas.

5. Conclusion

It is possible to build distributed file systems that use a single core protocol to provide both scale and notions of security. The Folklore distributed file system is built around scalable endemic replication, whose chaotic yet self-stabilizing behavior provides replica untraceability. Experiments show resilience to massive attacks. We leave details of the Folklore membership protocol, and PC cluster-based experiments to a later article. We are also augmenting Folklore to withstand other types of attacks.

References

- [1] *Blowfish*, <http://www.schneier.com/blowfish.html>.
- [2] *SHA-1*, <ftp://ftp.rfceditor.org/in-notes/rfc3174.txt>.
- [3] R. Bhagwan et al. Total Recall: System Support for Automated Availability Management. In *Proc. NSDI*, pages 337–350, 2004.
- [4] R. Bhagwan, S. Savage, and G. Voelker. Understanding availability. In *Proc. IPTPS '03*, Feb. 2003.
- [5] M. Castro and B. Liskov. Practical Byzantine fault tolerance. In *Proc. OSDI '99*, pages 173–186, 1999.
- [6] A. Das, I. Gupta, and A. Motivala. SWIM: Scalable Weakly-consistent Infection-style Process Group Membership Protocol. In *Proc. DSN*, pages 303–312, 2002.
- [7] J. R. Douceur. The Sybil Attack. In *Proc. IPTPS '01*, pages 251–260, 2002.
- [8] I. Gupta. On the design of distributed protocols from differential equations. In *Proc. ACM PODC*, pages 216–225. ACM Press, 2004.
- [9] J. Kubiawicz et al. OceanStore: an architecture for global-scale persistent storage. *SIGARCH Comput. Archit. News*, 28(5):190–201, 2000.
- [10] P. Maniatis et al. Preserving peer replicas by rate-limited sampled voting. In *Proc. SOSP '03*, pages 44–59, 2003.
- [11] D. S. Roselli, J. R. Lorch, and T. E. Anderson. A comparison of file system workloads. In *USENIX Annual Technical Conference, General Track*, pages 41–54, 2000.
- [12] A. Rowstron and P. Druschel. Storage management and caching in PAST, a large-scale, persistent peer-to-peer storage utility. In *Proc. SOSP*, pages 188–201, 2001.
- [13] Y. Saito et al. FAB: building distributed enterprise disk arrays from commodity components. *SIGPLAN Not.*, 39(11):48–58, 2004.