

# Insurable Storage Services: Creating a Marketplace for Long-Term Document Archival

Rahul Simha<sup>1</sup> and K. Gopinath<sup>2</sup>

<sup>1</sup> Department of Computer Science  
The George Washington University  
Washington, DC 20052, USA  
simha@gwu.edu

<sup>2</sup> Department of Computer Science and Automation  
Indian Institute of Science  
Bangalore, India  
gopi@csa.iisc.ernet.in

## Abstract

Digital storage is a key element not only of computing systems, but is now considered an essential component of the infrastructure of any modern organization. This need has co-evolved with the technology that has grown rapidly in recent years to provide low-cost high-capacity storage. At the same time, the storage needs of users have now become more sophisticated and diverse. Some users require very long-term preservation; others need high security; and still others ask for highly-reliable, distributed storage solutions. These needs pose a problem for solution providers in that no single solution seems to meet all needs. Similarly, users must construct services out of disk systems on their own. This paper proposes a way to streamline the marketplace through *insurable storage services*, a combination of two ideas. The first is to define different categories of storage service; the assumption here is that a refined categorization will better identify particular user needs. The second, and more substantive idea, is to treat digital documents as insurable property. The insurance of storage will provide economic incentives for both producers (storage service providers) and consumers (individuals, organizations) to jointly create a marketplace that provides a diversity of differentially-priced services. For example, insurers can help assess the durability of storage solutions and provide consumers with a quantitative valuation (“It’ll cost you \$x per GB to ensure that your documents last 100 years”). Similarly, storage service providers will have incentives to maintain multiple geographically distributed copies, and to continually move the copies onto emerging technologies (“You’ll need to store more copies if you want a higher reliability rating”).

## Keywords

Digital storage, storage networks, secure storage, disk systems, document archival

## 1 Introduction

Digital storage has grown from being a mere “part of a computer” to the complex array of interacting technologies, file systems, services and service providers that defines the world of storage today. This world now has evolved into sub-specialties such as materials, solid-state memories, disk-systems, RAID, storage area networks and storage services, to name a few [26].

In this paper, we identify some possibilities for future storage services, and most importantly, focus on a market mechanism by which these services could arise and grow naturally. The lack of such a dynamic marketplace is one of several key problems identified in the recent report on digital preservation [18]. In the words of that report, “Creating an economy for long-term preservation entails providing incentives for organizations to invest in digital archives, even though some of the benefits of investments made today may not be realized for decades.”

Our approach to addressing this problem starts with the observation that most current storage-service paradigms posit two players in the service market: the user and the Storage Service Provider (SSP). By introducing a third player, insurers, and by treating digital documents as property, we argue that a natural marketplace will emerge as the insured and the insurers drive the creation of new, cost-effective services. This way, investments can be funded today for the benefits realized later. The purpose of this paper is to provide an overview of this idea, to identify factors affecting insurance and to raise a few new technical issues.

## 2 Background

### 2.1 System components

As background, let us outline the various components and players in a storage service system. Consider a typical desktop user who is busy creating, updating and deleting files. This creates a stream of I/O between processor and disk. For many users, the local disk is the final resting place for documents instead of being merely a local cache of a more stable repository. In today’s computing paradigm, individual users either perform their own backups or are part of a networked backup system. In either case, files are scheduled for backup on a periodic basis. When the individual user is part of a larger organization, the organization may contract with a SSP to perform such a backup, or even store the primary copy on site.

In our paradigm, we identify the following players. First, there are *users*, either as individuals or as part of an organization. Each user machine runs software to feed files into a pipe intended for storage. That is, the operating system, instead of writing to a local disk hands over the file to an application software running on behalf of a service provider. This application will use the local disk as a cache for efficiency. Second, *SSP’s* at the other end of this pipe provide a variety of storage services, including storing the primary file and several copies.

Finally, some services will need third-party providers. These include: trusted third-parties for some types of security operations (we will see one example later in Section 4.1), third-party services for search, or a third-party curation service aimed at long-term archival for maintaining a historical record.

## 2.2 User needs

What do users want out of a storage provider? Let us first consider individual users in the home or small-business environment:

- *Backup.* Clearly, this most common need today will continue into the future. In addition to simple backup, different users may desire various levels of quality. For example, small business users might opt for multiple, geographically diverse backup copies.
- *Legacy.* Users will want to pass on their digital property to rightful heirs and ensure their access. Some of these heirs may be family members, others might include a trust or a public organization. Today, this type of service is virtually non-existent and certainly not systematically addressed.
- *Legal services.* As with any kind of property, in case of conflict over a legal document, storage providers should make available a framework for conflict resolution between all those party to the document. For example, a real-estate property deed will need to have its integrity assured, and if needed, to be checked against the same document stored by the local county office<sup>1</sup>. Thus, a SSP must provide the a means for access by other interested parties, including a court of law.

In addition to these services, organizations have other requirements:

- *Strong backup and availability.* Organizations usually place greater value on the accessibility and reliability of backups than do individuals. As noted by others [7, 17], these additional factors are influenced by the physical durability of storage warehouses, their geographical diversity, diversity across disk manufacturers, number of such diverse copies, administrative convenience and speed of archival.
- *Group coordination.* Organizations also need to store documents in collections and from groups of individuals, often removing association with individuals. For this purpose, an important service is to provide organizations with a framework for defining groups and moving documents between groups. Organizations also share documents between divisions and across other organizations.
- *Ownership.* The question of ownership is important to an organization. When documents are associated with groups or roles, individuals in those groups or roles assume temporary ownership. Thus, a framework for ownership must provide for associations between individuals and groups, and between individuals and administrative roles with high privileges.

---

<sup>1</sup> Although cryptographic techniques can be used to verify integrity, these techniques are subject to implementation errors. Furthermore, challenges must nonetheless be addressed by third parties.

The needs of public organizations or of government agencies also include custodial actions. For example, U.S. law dictates that the executive branch must make available its digital record to the National Archives. Even if the law does not require it, other agencies might wish to make their documents available to the public after a period of use.

## 2.3 Services

Each need identified above corresponds to a service that an SSP should provide. In addition to these, we identify a few more:

- *Security.* SSP's should go beyond simple encryption to provide a comprehensive solution to various security needs. These include key management, key transfers between groups, long-term key storage, re-keying, integrity checking, audit trails, conflict resolution and a framework for access privileges in a multi-user environment.
- *Escrow.* SSP's should provide services that aim for fixed-term secure storage, after which documents become public or are made available to other parties according to a policy. Escrow services could also be used in various legal transactions, such as on-line contract negotiations.
- *Notification services.* SSP's should provide periodic notification, just like a bank statement, that contain storage transaction histories, reports on integrity checks, updates of technologies, numbers of copies, the geographic sites where they are stored and usage statistics.
- *Search.* SSP's should provide access points for search providers so that sophisticated search technologies may be used by organizations to search within their document troves.
- *Meta information.* SSP's should provide API's for operating systems and applications so that they may store context information along with documents. This context information should be augmented with usage statistics, locality statistics (time, date) to form a comprehensive record of meta information that will be useful for later searches.

Taken together, these features point to a future of value-added services in the area of storage. Why hasn't this future emerged already? Every change in the marketplace needs a driver, a force that facilitates transactions and growth and propels providers to competitively offer creative new services. We argue that one such possible market driver is insurance.

## 3 Insured Storage

### 3.1 How it would work

Users and organizations today are faced with a bewildering number of options in constructing solutions for their storage needs. They must identify vendors of storage systems and then build services on top of these storage systems. Furthermore, they have no easy way to quantify the value of their investment.

Instead, we propose that users treat their digital documents as property, whose storage is insurable. To illustrate, let us consider a user interested in contracting out storage services from SSP's. The user obtains an insurance rating for various service options across vendors, and together with the prices for each, makes a decision. Thus, for example, Insurance Company *A* is willing to insure service *B* offered by SSP *C* at the rate of \$10 per MB for 5 years. Faced with choices such as these, users can discern what works best for them.

Similarly, SSP's competing in this marketplace will make optimal use of technologies to acquire the best possible insurance rating and to offer differentiated services at different prices. Thus, SSP *C* can offer a long-term archival service (50 years) but with slow access speeds at \$10 per MB and a short-term (5 years), highly efficient service for \$40 per MB. Each of these might be rated differently by insurers. To get a high insurance rating, a SSP would have to convince insurers that, for example, they are using current technologies, exploiting geographic and vendor diversity, and creating fresh copies at reasonable time intervals.

### 3.2 Benefits of insurance

The chief benefit of using insurance is that it mediates between users and providers in what could be a confusing array of options and benefits. Users are given a single number (or two) by which to assess services; the trustworthiness of this number is based on the reputation of insurers and the cost of insurance. Similarly, SSP's work towards better ratings by cycling copies onto new technologies, increasing the number of copies, using different hard-drive vendors, different geographic locations for the copies, physically securing these locations from natural catastrophes and using proven engineering practices. All of these can be fitted into statistical models that quantify reliability, availability and efficiency [9]. Insurers, with their legions of statisticians and considerable experience with reliability models, are already suited to this type of analysis. Furthermore, SSP's will be driven to providing creative new services to meet the needs of customers, as and when these get included in the overall rating.

We conjecture that an added benefit for long-term archival would be some closure to the "format wars." As is well-known in the archival literature, there is a tension between the proprietary formats (such as Microsoft Word) that users prefer to use, and the open, human-readable formats (such as XML) that, according to archivists, offer the best hope of being readable in the future [5, 16, 21, 23]. Since a proprietary format should logically receive a low 100-year rating, users will be driven towards open formats if they are to guarantee long-term storage.

Note that other mechanisms exist that could mediate between consumers and producers in the storage marketplace, as can be found in other marketplaces. For example, a softer, non-binding form of mediation is to simply provide ratings, of the form provided by consumer protection groups. However, these are not as powerful as mechanisms such as insurance that associate high cost for inaccurately rating services.

### 3.3 Factors affecting insurance

Many factors that one should use in rating the quality of a storage service have been discussed in other papers [7, 8, 17]. These include the obvious ones: number of copies, geographic diversity of copies, vendor diversity, hardware reliability, copy replacement policy, hardware replacement policy and physical security. To these we add a few new ones. These are all based on the notion that a document that can't be deciphered or precisely located in the future, even if it is known to exist, is useless.

- *Formats.* As mentioned earlier, proprietary formats, especially those subject to frequent change, are probably not suited for long-term storage. Even if these must be used, proper emulation [16, 21, 22] can mitigate some of the disadvantages, and thus should be factored into a rating. Even among open-formats, documents that use a mixture (such as a combination of HTML and JPEG) might receive a lower rating than simple documents.
- *Searchability.* With a long-term view, documents should be retrievable not just by specifying the owner but through a variety of means. A custodian of government documents, for example, should be certain that the documents will show up in a variety of search approaches.
- *Monitoring quality.* SSP's that continuously monitor their document copies for integrity should obviously receive a better rating. Monitoring should also include hardware, network connections, software patches and security in its purview.
- *Networks.* Aside from the assessment of individual sites where copies are stored, the network used to connect these sites should also be subject to evaluation. SSP's that use multiple networks with redundant connections should get a higher rating.

## 4 New technical issues

There are a host of technical issues associated with storage systems and storage services. Many of these are taken up in papers on individual issues or have been surveyed elsewhere [1, 2, 7, 11, 19, 20, 23–25]. Here, we point out two technical issues related to long-term archival that we believe are relatively new and not yet addressed in the technical literature.

### 4.1 Ownership devolution

The first issue is motivated by considering those situations in which ownership of a document must be transferred. For example, an individual user might leave her digital assets to heirs upon her demise. Similarly, the executive branch of government must leave their documents to the public after their term expires [6]. The problem is non-trivial even when these documents are not encrypted during use – there must be a way, for example, to ensure that the policy is properly enforced. In the future, however, most documents will probably be encrypted

when stored using a service. In this case, how can one guarantee that the keys will be made available? What happens if the keys are lost or if an individual loses them?

To address these issues, SSP's need a framework for handling ownership devolution. To solve problems such as key escrow, trusted third parties might be needed in addition to SSP's. For example, a trustee can perform encryption enroute to storing a document; this trustee can keep the key and be required to reveal the key at some point in the future. A trust management framework might therefore become an essential part of such services. Since access control is based on identity, third party input is needed so that trust, delegation and public keys can be negotiated.

Note that with public-key cryptography, it becomes possible to deal with anonymous users as long as they have a public key: authentication and authorization are now possible with models such as SDSI/SPKI[10]. In this case, an issuer authorizes specific permissions to specific principals; these credentials can be signed by the issuer to avoid tampering. For example, SDSI/SPKI provides for credentials with delegation with the assumption that locally generated public keys do not collide with other locally generated public keys elsewhere. This allows exploiting "local namespaces": any local resource controlled by a principal can be given access permissions to others by signing this grant of permission using the public key.

Access control and cryptography can now be combined into a larger framework with logic for authentication/authorization and access control. For example, an authorization certificate  $(K, S, D, T, V)$  in SDSI/SPKI can be viewed as an ACL entry, where keys or principals represented by the subject  $S$  are given permission, by a principal with public key  $K$ , to access a "local" resource  $T$  in the domain of the principal with public key  $K$ . Here,  $T$  is the set of authorizations (operations permitted on  $T$ ),  $D$  is the delegation control: whether  $S$  can in turn give permissions to others and  $V$  is the duration during which the certificate is valid. Name certificates define the names available in an issuer's local name space whereas authorization certificates grant authorizations, or delegate the ability to grant authorizations. A certificate chain provides proof that a client's public key is one of the keys that has been authorized to access a given resource either directly or transitively, via one or more name-definition or authorization-delegation steps.

However, permission-based trust management cannot authorize principals with a certain property easily. For example [14], to give a free digital copy of a book to students, the university bookstore can delegate "free copy" permission to the institute key. The university has to delegate its key to each student with respect to "bookstore" context; this can place too high a burden on the institute. The cost is likewise high if the institute creates a new group key for students and delegates it to each student key. One solution to this problem is attribute-based approach: it combines RBAC and trust management. Other approaches include proof carrying authentication [3].

## 4.2 Universal document ID

Currently, documents are identified by user account, directory structure and file name. However, as storage services and providers grow, users will need to transfer documents across providers and systems. Furthermore, there are many situations where documents need to be identified with groups rather than individuals. This raises the issue of proper identification of documents for long-term use. Identifiers are needed for indexing, for maintaining meta information and for tracking.

Thus, a service that provides unique universal document identifiers will be useful. At the same time, these identifiers should maintain an individual's privacy and not be traceable back to the individual. One solution is the use of SHA2 hashes such as sha256, sha384 and sha512 algorithms. If SHA2 is strongly collision-resistant (as is currently believed), then the hash can be used as a global identifier. Usability of such hashes can be managed with a secure mapping between names and hashes that is locally maintained. The latter service could be constructed using a mechanism similar to DNS [15], the mechanism for internet domain names. If HMAC-SHA2 is used with the secret being between the client and the storage provider service, both origination and integrity can also be guaranteed.

One can imagine further refinements. These IDs, if properly extended, can also be used as ways of authorizing the use of documents to others (as in DRM). One possibility is the use of "capabilities." However, undesirable information flows are possible in such systems using capabilities in the presence of Trojans [12]. This requires that some information about the intended recipient is incorporated in the capability, and thus results in a modification of the strict capability model.

For example[4], assume that there are clients, servers and metadata servers in a system. Metadata servers provide information about the information present in the servers. Assume there are secure channels between client and metadata server and between servers and client (using symmetric keys). A capability will have the object ID and the permissions (such as read, write or archive). A credential sent to the client will have the capability along with the encryption of a secret  $K$  by the secret key established between server and metadata server, and validation tags that are based on MAC.  $K$  is sent to the client on the secure channel. This secret  $K$  essentially makes the capability private to the specific client. The credential is sent by the client to the server along with a signed MAC (HMAC-SHA2) on the channel name using the secret  $K$  as key. Since only the client has  $K$  the server can be sure of the credential sent by the client. To go from the HMAC-SHA2 to the document, the server has to keep a content addressable map from hashes to documents as happens in many cryptographic filesystems but this can be at whole-file level rather than at the block level. If SHA2 is used instead of HMAC-SHA2, the document can be public with only integrity being the issue.



## 5 Summary

In this paper, we have described the benefits of using insured storage as a mechanism for driving the marketplace for storage services. In addition, we have identified a few new factors affecting insurability and outlined a couple of new technical issues relating to long-term storage.

The question that naturally arises is: will users take to insuring their documents? Our view is that, except for a few very important documents, users might be unwilling to negotiate such minutiae on a per-document basis. However, we speculate that users will be willing to insure services especially if offered on an annual basis or in terms of Gigabytes, similar to some internet services. Ultimately, the services will be offered (because they are needed) and ultimately, users will find some, perhaps informal, way of assessing their value. Insurability, in addition to opening up a new possibilities for that industry, offers a way to rapidly streamline the marketplace for storage services. Already, insurance is currently being offered for costs related to compliance with open-source software standards [13].

Finally, we recognize that associating a direct cost with storage might drive users towards placing a value on documents today, resulting in the loss of some documents that could be valuable to some future historian. This is an important issue in general that is not considered in this paper.

## References

1. A.Adya, W.J.Bolosky, M.Castro, G.Cermak, R.Chaiken, J.R.Douceur, J.Howell, J.R.Lorch, M.Theimer and R.P.Wattenhofer. FARSITE: Federated, available, and reliable storage for an incompletely trusted environment. *Proc. 5th Symp. Operating Systems Design and Implementation*, Boston, MA, December 2002.
2. G.A.Alvarez et al. Minerva: an automated resource provisioning tool for large-scale storage systems. *ACM Trans. Comp. Systems*, 2001.
3. A.W.Appel and E.W.Felten. Proof-carrying authentication. *6th ACM Conference on Computer and Communications Security*, November 1999.
4. Alain Azagury, Ran Canetti, Michael Factor, Shai Halevi, Ealan Henis, Dalit Naor, Noam Rinetzky, Ohad Rodeh, and Julian Satran. A Two Layered Approach for Securing an object store network. *1st International IEEE Security in Storage Workshop (SISW 2002)*.
5. D.Bearman. Reality and chimeras in the preservation of electronic records. *DLIB Magazine*, Vol.5, No.4, 1999.
6. D.Bearman. The implications of Armstrong vs. Executive Office of the President for the archival management of electronic records. *American Archivist*, 56, 1993, p.150-160.
7. B.Cooper, A.Crespo and H.Garcia-Molina. Implementing a reliable digital object archive. *Proc. Fourth European Conference on Digital Libraries (ECDL)*, Sept, 2000.
8. A.Crespo, H.Garcia-Molina. Cost-driven design for archival repositories. *Proc. 1st Joint Conference on Digital Libraries (JCDL)*. June, 2001.
9. A.Crespo and H.Garcia-Molina. Modeling archival repositories for digital libraries. *Proc. Fourth European Conference on Digital Libraries (ECDL)*, Sept, 2000.

10. C.Ellison, B.Frantz, B.Lampson, R.Rivest, B.Thomas and T.Ylonen. SPKI Certificate Theory. Internet Network Working Group RFC2693, September 1999.
11. G.R.Ganger, P.K.Khosla, M.Bakkaloglu, M.Bigrigg, G.R.Goodson, S.Oguz, V.Pandurangan, C.Soules, J.D.Strunk and J.J.Wylie. Survivable storage systems. DARPA Information Survivability Conference and Exposition (Anaheim, CA, 12-14 June 2001), pages 184-195 vol 2. IEEE, 2001.
12. S.Halevi, P.Karger and D.Naor. Enforcing Confinement in Distributed Storage and a cryptographic model for access control, *IBM Tech Report*, 2002.
13. M.LaMonica. Insurer launches \$10 million open-source policy. CNET News, Oct 2005. [http://news.zdnet.com/2100-3513\\_22-5924112.html](http://news.zdnet.com/2100-3513_22-5924112.html)
14. N.Li, B.Grosz and J.Feigenbaum. Delegation Logic: a logic-based approach to distributed authorization, *ACM Trans. on Info. and System Security*, Feb 2003.
15. C.Liu, P.Albitz and M.Loukides. DNS and BIND. *O'Reilly Pub.*, 1998.
16. R.A.Lorie. Long-term archiving of digital information. *IBM Research Report RJ-10185*, May 2000.
17. R.W.Moore, J.F.Jaja and R.Chadduck. Mitigating risk of data loss in preservation environments. Prof. 22nd NASA Conf. on Mass Storage Systems and Technologies (MSST), 2005.
18. It's About Time: Research Challenges in Digital Archiving and Long-term Preservation National Digital Information Infrastructure and Preservation Program ([www.digitalpreservation.gov](http://www.digitalpreservation.gov)).
19. Technical Architecture, Version 0.2 National Digital Information Infrastructure and Preservation Program ([www.digitalpreservation.gov](http://www.digitalpreservation.gov)).
20. J.Park and R.Sandhu. The UCON-ABC usage control model. *ACM Trans. Info. Sys. Security*, Vol.7, No.1, Feb 2004, pp.128-174.
21. J.Rothenberg. Ensuring the longevity of digital documents. *Scientific American*, Vol. 272(1), January 1995.
22. J.Rothenburg. An experiment in using emulation to preserve digital publications. Technical Report, *RAND-Europe*, April 2000.
23. M.Smith. Eternal bits. *IEEE Spectrum*, July 2005, pp.16-21.
24. V.Sriram. SAFIUS: A secure accountable filesystem for untrusted storage. MS Thesis, Indian Institute of Science, Bangalore, 2004.
25. Q.Xin, E.L.Miller, T.Schwarz, D.E.Long, S.A.Brandt and W.Litwin. Reliability mechanisms for very large storage systems. Proceedings of the 20th IEEE / 11th NASA Goddard Conference on Mass Storage Systems and Technologies (MSST 2003), San Diego, CA, April 2003, pages 146-156.
26. J.Tate, R.Kanth and A.Telles. An introduction to storage area networks. IBM Redbook. [www.ibm.com/redbooks/](http://www.ibm.com/redbooks/)