

Project no. 269977

**APARSEN**  
**Alliance for Permanent Access to the Records of Science**  
**Network**

**Instrument:** Network of Excellence

**Thematic Priority:** ICT 6-4.1 – Digital Libraries and Digital Preservation

**D27.1 RECOMMENDATIONS ABOUT  
SCALABILITY**

---

Document identifier:	<b>APARSEN-REP-D27_1-01-0_6</b>
Due Date:	2013-07-31
Submission Date:	2013-08-07
Work package:	WP27
Partners:	CSC, BL, TESSELA, IBM
WP Lead Partner:	IBM
Document status	Released
URN	urn:nbn:de:101-20140516208

---

**Abstract:** This report presents the results of the analysis of scalability as part of APARSEN WP27 research. The report gives the background on scalability; uses survey results to identify scalability needs and gaps; and provides recommendations on scalability of tools and services.

**Delivery Type** Report  
**Author(s)** CSC, BL, TESSELA, IBM  
  
**Approval** David Giaretta, Simon Lambert (coordinators)  
**Keyword List** Digital preservation, Storage, Scalability, Cloud  
**Availability** ☒ PUBLIC

#### Document Status Sheet

Issue	Date	Comment	Author
0.1	2013-07-03	Initial version by IBM	O. Edelstein
0.2	2013-07-22	Added more sections	O. Edelstein, A. Hanter, D. Naor, S Cohen, M Factor, D. Harnik
0.3	2013-07-29	Editing, more data	O. Edelstein
0.4	2013-08-01	Version for review	O. Edelstein
0.5	2013-08-02	Internal review	K.Kaur
0.6	2013-08-06	Final with review comments addressed	O. Edelstein

---

**Project information**

Project acronym:	<b>APARSEN</b>
Project full title:	<b>Alliance for Permanent Access to the Records of Science Network</b>
Proposal/Contract no.:	<b>269977</b>

---

**Project coordinator: Simon Lambert/David Giaretta**

Address:	STFC, Rutherford Appleton Laboratory Chilton, Didcot, Oxon OX11 0QX, UK
Phone:	+44 1235 446235
Fax:	+44 1235 446362
Mobile:	+44 (0) 7770326304
E-mail:	<a href="mailto:simon.lambert@stfc.ac.uk">simon.lambert@stfc.ac.uk</a> / <a href="mailto:david.giaretta@stfc.ac.uk">david.giaretta@stfc.ac.uk</a>

## CONTENT

<b>EXECUTIVE SUMMARY .....</b>	<b>5</b>
<b>1 INTRODUCTION .....</b>	<b>6</b>
1.1 DOCUMENT STRUCTURE .....	6
<b>2 BACKGROUND .....</b>	<b>7</b>
2.1 SCALABILITY REQUIREMENTS .....	7
2.2 SCALABILITY IN SELECTED PRESERVATION SYSTEMS AND TECHNOLOGIES .....	11
2.2.1 Commercial preservation products and services .....	11
2.2.1.1 Tessella's Safety Deposit Box & Preservica SaaS solutions .....	12
2.2.1.2 Ex Libris Rosetta .....	13
2.2.2 Open Source .....	15
2.2.2.1 Duraspace .....	15
2.2.3 EU Research Projects .....	16
2.2.3.1 ENSURE .....	16
2.2.3.2 SCAPE .....	16
2.2.3.3 SCIDIP-ES .....	17
2.3 CLOUD STORAGE .....	17
2.3.1 What is Cloud Storage? .....	17
2.3.2 Scalability in Cloud Storage .....	18
2.3.2.1 Amazon S3 .....	18
2.3.2.2 CloudFiles Rackspace .....	19
2.3.2.3 Google Cloud Storage .....	19
2.3.2.4 iCloud .....	20
2.3.2.5 DropBox .....	20
2.3.2.6 OpenStack Swift .....	20
2.3.3 Cloud Delivery Models .....	20
<b>3 SCALABILITY ANALYSIS .....</b>	<b>22</b>
3.1 GENERAL REMARKS .....	22
3.2 SCALABILITY PARAMETERS .....	22
3.3 RESPONDING ORGANIZATIONS .....	23
3.4 RESULTS AND ANALYSIS .....	23
3.4.1 Data Types Preserved .....	23
3.4.2 Size of the Data .....	25
3.4.3 Number of Digital Objects and Their Sizes .....	27
3.4.4 Versions .....	28
3.4.5 Connections between Objects .....	29
3.4.6 Size of the Metadata and Connections between Metadata .....	29
3.4.7 Retention Period .....	30
3.4.8 Availability (Access Rate, Access Location and Access time) .....	31
3.4.9 Distribution .....	32
3.4.10 Searchability .....	33
3.4.11 Degree of sharing .....	33
3.4.12 Overall Results Analysis .....	33
<b>4 RECOMMENDATIONS .....</b>	<b>35</b>
<b>5 ANNEXES .....</b>	<b>37</b>
5.1 QUESTIONNAIRE .....	37
5.2 ANNEX II: QUESTIONNAIRE INSTRUCTIONS .....	46
5.3 ANNEX III: QUESTIONNAIRE GLOSSARY .....	51
<b>6 REFERENCES .....</b>	<b>53</b>

## EXECUTIVE SUMMARY

The amount of data produced in the world is growing rapidly. Organizations are facing the need to preserve more and more data. The complexity of the data to preserve is increasing as well – more types of data and more connections among the data. Additionally, the requirements from the preservation system are more demanding. In many cases the data should be available at any time, in any place, and immediately. The high increase in scalability demands pose huge challenges to the current preservation systems.

This report presents the work done on understanding scalability demands for the next decade and suggesting where attention should be focused to help to cope with the challenges. This work complements the work done in APARSEN in other areas such as tools, services, and storage solutions, concentrating on the overall scalability aspects of preservation systems.

The first step in understanding scalability was to identify aspects that play a role in scalability. A list of the scalability related parameters was defined:

- Variety of data types
- Size of the data
- Number of digital objects and size of each object
- Number of versions of the same object
- Connections between different objects
- Amount of metadata and connections between metadata
- Retention Period
- Availability
- Distribution
- Searchability
- Degree of sharing

Scalability requirements of organizations were collected from existing systems using surveys. The results were analyzed to identify gaps. State of the art scalable solutions such as cloud storage were contrasted to the existing state of the surveyed systems. Several recommendations are suggested:

- Scalability should be planned as an integral part of the tools and services development.
- Whenever possible, the use of commercially available systems, tools, and services should be considered instead of using in house developed solutions
- In order to benefit from economies of scale, the use of cloud storage should be considered
- Consider preservation as part of the overall organization's activity.

## 1 INTRODUCTION

This report summarizes the work on scalability done as part of APARSEN WP27. This is part of the usability topic in APARSEN covering activities that relate to the usability of preservation systems: standards, tools, interoperability, and scalability.

The objective of the work is to understand scalability gaps in long term digital preservation technologies that need to be addressed in order to attain the required levels of scalability. The report summarizes the challenges from, and possible responses to, scalability needs over the next decade.

When looking at scalability, we took a broad view. In addition to the obvious factor of size of the data, we considered the complexity of the data such as the interconnections among data and the variety in data formats and sizes, as well as the requirements from preservation system in terms of availability and accessibility.

The main steps of the work were:

- Understand what the important scalability parameters are in preservation systems
- Understand the scalability requirements of the preservation systems for the next few years
- Identify gaps in technology that prevent us from getting to the right level of scalability
- Summarize challenges and recommend areas that need to be addressed.

The work on scalability is important to enable better support for the rapid increase in the amount of data that need to be preserved and its growing complexity. The recommendations report will be used by the VCoE to direct activities towards covering the gaps and requirements identified.

### 1.1 DOCUMENT STRUCTURE

The structure of this document is as follows:

Chapter 2 gives background on scalability. It presents general information on scalability requirements, the scalability aspect of several systems and research projects, and scalability in cloud systems.

Chapter 3 gives the results of a survey used to understand scalability requirements and status and analysis of the results.

Chapter 4 gives our recommendations on scalability of preservation systems and services.

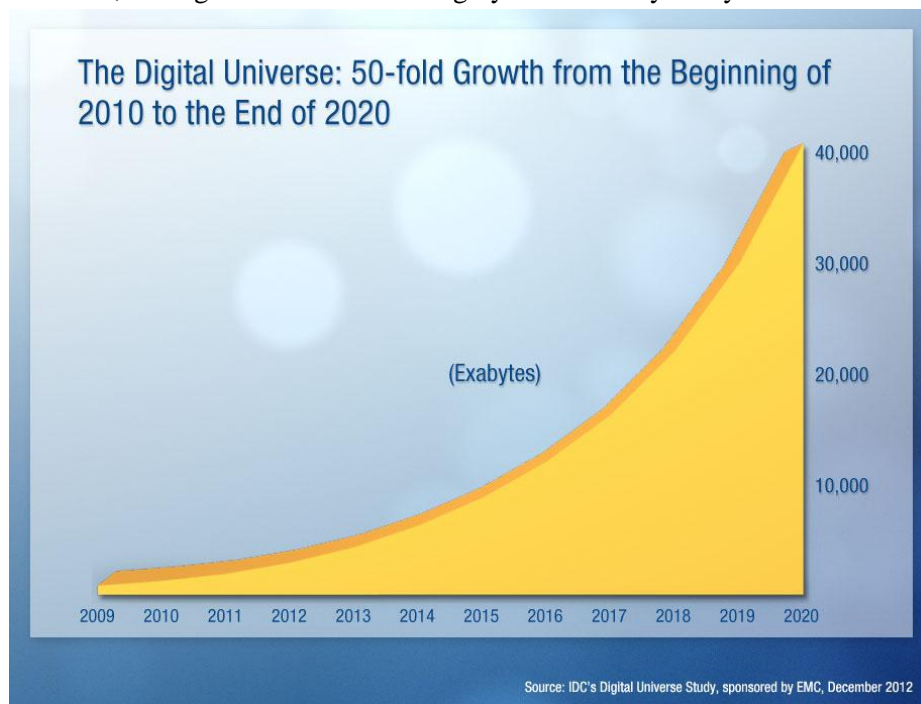
The annexes detail the questionnaire used for the survey and the relevant terminology.

## 2 BACKGROUND

### 2.1 SCALABILITY REQUIREMENTS

This section presents information on the expected spiraling growth of data for the coming decade. The predicted growth poses a huge challenge for the scalability of preservation systems.

IDC published a series of studies on the growth of digital data across the world. The first study was in 2007, the second in 2010, and the latest in 2012 [3] [4] [5]. The 2012 study gives predictions on the data in 2020 and shows that from 2005 to 2020, the digital universe will grow by a factor of 300, from 130 Exabytes<sup>1</sup> to 40 Zettabytes (more than 5,200 gigabytes for every man, woman, and child in 2020). From now until 2020, the digital universe will roughly double every two years:

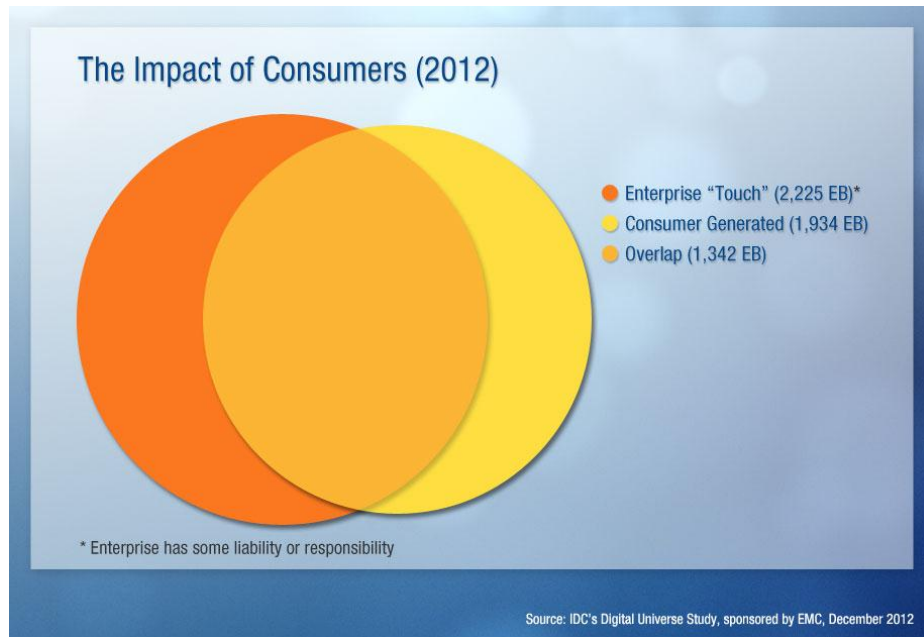


Whilst this forecast may look high, looking back at prior forecasts, the reality is that the actual growth is higher than the predicted growth. The 2007 study predicted less than a Zettabyte of digital data in the world in 2010, while the actual was 1.2 Zettabytes. Predictions are getting higher too; in 2010 IDC predicted 35 Zettabytes of data in 2020, in the 2012 study the forecast was increased to 40 Zettabytes.

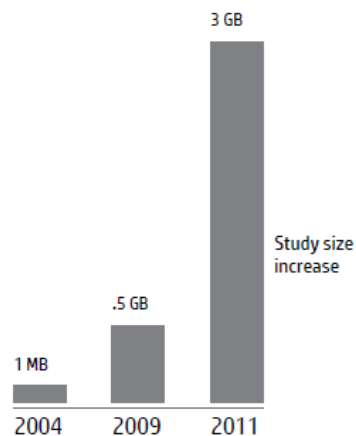
The data ranges from images and videos to digital movies populating the pixels of high-definition TVs to transponders recording highway tolls. It also includes the more traditional corporate data, such as banking, healthcare, and scientific data.

Most of the digital information is created by consumers; however, enterprises still have liability for the vast majority of the information (indicated in the chart below) as they have to deal with the issues of copyright, privacy, and compliance.

<sup>1</sup> An Exabyte is  $2^{60}$  Bytes; a Zettabyte is  $2^{70}$  Bytes or 1024 Exabytes.



While the preservation of large parts of digital information is not considered vital, the information that does need to be preserved is also growing rapidly. One area is in the healthcare domain. The size and quantity of images is growing fast. Those images need to be preserved according to medical practices and legal regulations for at least several years. The two charts below give some data on the growth.



**Medical images will require  
30% of the world's storage  
(Source: HIMSS Analytics).**



## Current Digital Archiving Growth Problems

### o Digital Data Growth:

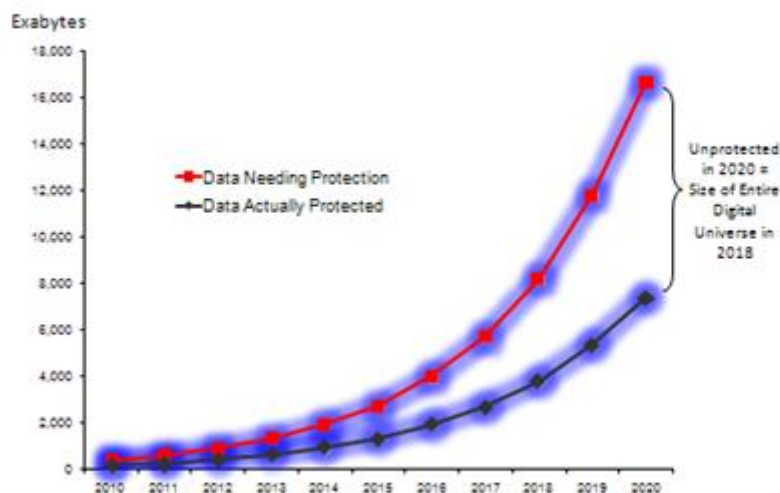
#### o A Radiologist Read. . .

o 1975 CT scans 80 images per day	0.4 MB
o 1984 CT scans 704 images per day	45.12 MB
o 1998 CT scans 1,600 images per day	993.6 MB
o 1998 MR scans, 1,314 images per day	102.6 MB
o 2006 CT scans, 4,650 images per day	595.2 MB
o 2006 MR scans, 11,800 images per day	4,530.0 MB
o 2009 CT scans, 22,640 images per day	2,677.2 MB
o 2009 MR scans, 11,830 images per day	6,065.7 MB
o 2009 PT scans, 2,127 images per day	2,563.5 MB

The requirements for availability of the data are also increasing. According to the IDC study, more and more organizations require that data will be produced in response to requests from any remote site in an acceptable amount of time. This poses scalability requirements on the systems holding the data.

Another requirement is data protection. The proportion of data in the digital universe that requires protection is growing at a faster rate than the total amount of data, from less than a third in 2010 to more than 40% in 2020. The chart below shows this trend and the growing gap.

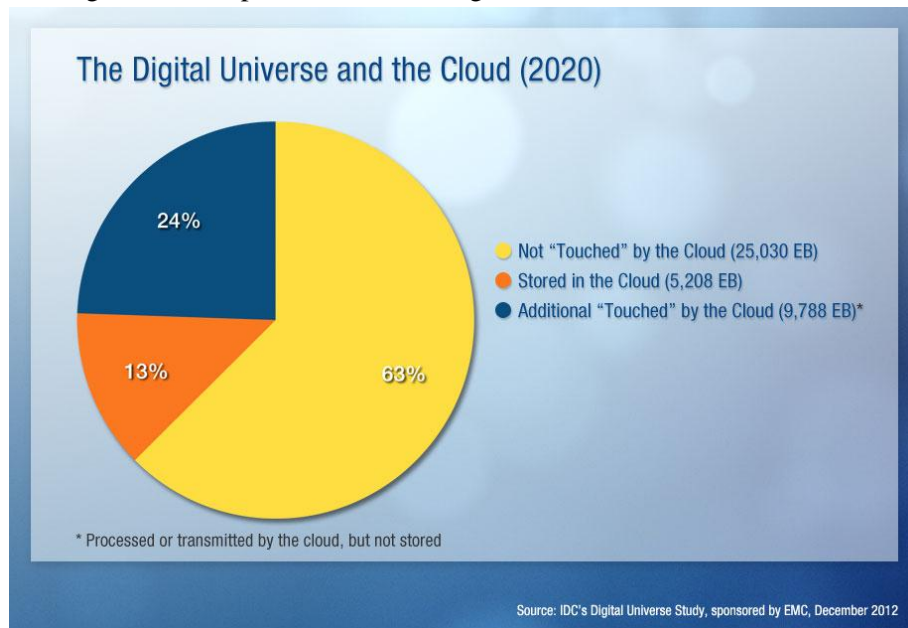
## Unprotected Data Needing Protection



Source: IDC Digital Universe Study, sponsored by EMC, May 2010

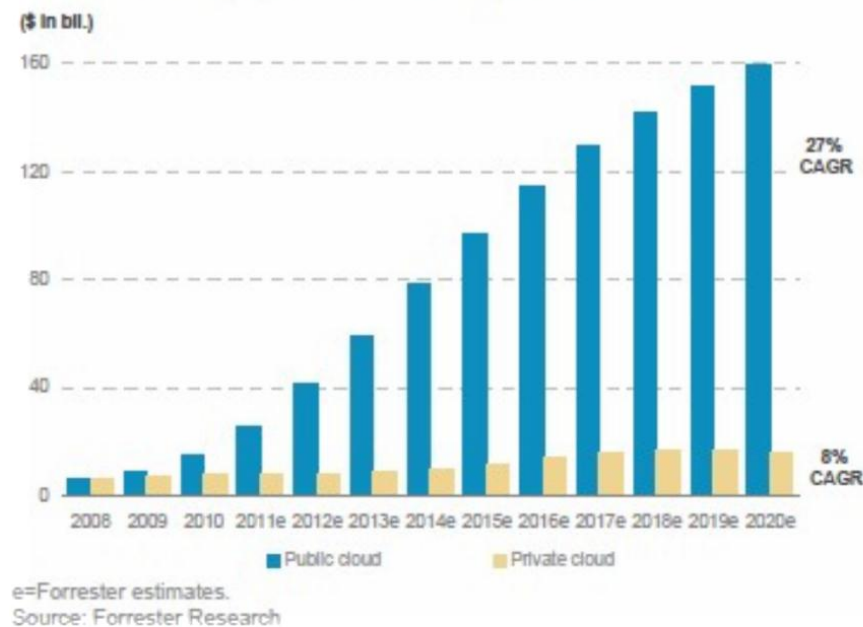
Cloud storage is becoming more pervasive. Public cloud growth rate is exponential, mostly for unstructured data, for example, pictures, documents, web pages, videos etc.

The IDC study predicts that by 2020, nearly 40% of the information in the digital universe will be "touched" by cloud providers — meaning that it will be stored or processed in a cloud somewhere in its journey from originator to disposal. The following chart illustrates this:

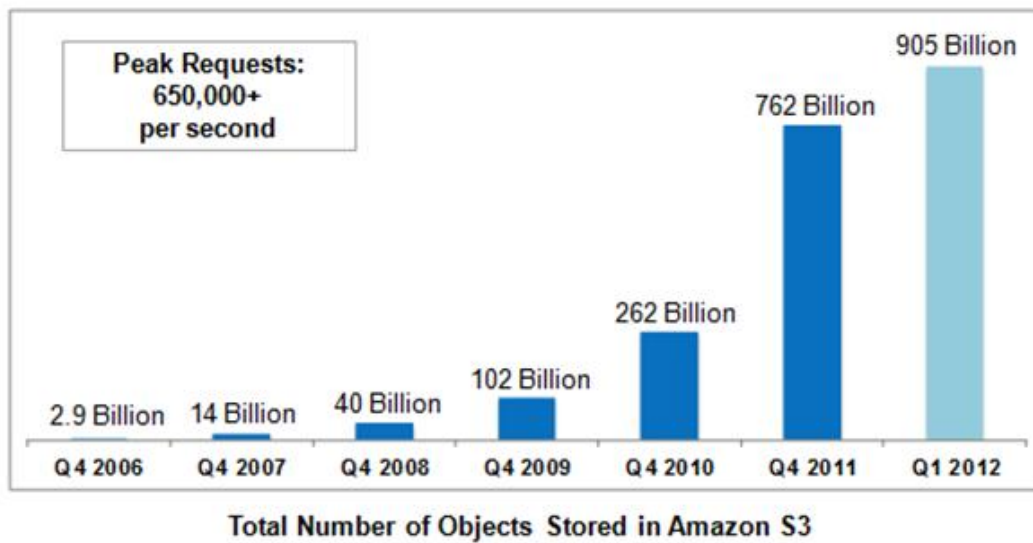


The expected growth of public cloud revenue is illustrated in the chart below.

### Forrester Expects Public Cloud (SaaS/PaaS/IaaS) to Reach Nearly \$160 Billion by 2020



The information from public cloud provider supports this trend. An Amazon report from April 2013 indicates that Amazon's S3 cloud storage now stores 2 trillion objects, up from 1 trillion in June 2012 and regularly peaks at over 1.1 million requests per second [6]. The chart presented below shows the growth in Amazon S3 stored data in the last few years.



The yearly statistics published by Facebook in March 2013 also indicate the growing amounts of data produced:

- Total users 1.1 Billion
- Daily active users 665 million
- Posts per month 1.8 Billion
- Total number of songs, albums or radio stations shared on Facebook: 110 million
- Total number of uploaded Facebook photos: 240 billion
- Average daily uploaded Facebook photos: 350 million

The spiralling growth in data and in requirements from the systems holding data, such as availability and data protection, impose a huge challenge on any system holding the data, including preservation systems.

## 2.2 SCALABILITY IN SELECTED PRESERVATION SYSTEMS AND TECHNOLOGIES

As the demand on archival systems increases due to the ever growing numbers of sources of digital information (data, documents, websites) and their corresponding volumes of content increase, the challenge of providing an IT based software solution and infrastructure to process and manage that massive volume of digital content within reasonable operational timescales becomes ever more apparent.

Content growth is not the only reason to be concerned with “operating at scale”, as user access requests and routine maintenance workflows can also drive up demand on the precious server side resources used to process those requests within the system. Preservation systems should be performing periodic background integrity checks on the stored content as well as its associated metadata, so a repository with 0.5 PetaBytes of digital content in it may consume significant CPU cycles, provided by dedicated servers, and a large slice of network bandwidth just to perform this vital task alone.

Providers of products and services have had to adapt the provision of archival systems and preservation services to meet the new frontier that is “the digital deluge”.

These products and services fall into three main camps: commercial products and service offerings; Open source offerings; and EC funded research projects. We now look at each of these in turn.

### 2.2.1 Commercial preservation products and services

There are a number of emerging commercial vendors that have developed digital preservation products

and services to specifically address the issues that arise when handling large volumes of data. These ‘scalable’ systems are typically modular in their design to facilitate parallelization wherever possible between their constituent component parts so that workflows are run over multiple server nodes, a process termed “Scale Out”. In addition to this, larger IT infrastructure can be deployed, with greater multiples of CPUs on the one server so as to improve the overall throughput of one particular service or application, a process known as “Scale Up”.

Combinations of “Scale Up” and “Scale Out” can be used to achieve the optimal performance for a specific archival or preservation system in response to anticipated usage requests and loadings.

We now consider some of the scalable digital preservation solutions that are currently available on the market that have been designed with scalability in mind.

### **2.2.1.1 Tessella’s Safety Deposit Box & Preservica SaaS solutions<sup>2</sup>**

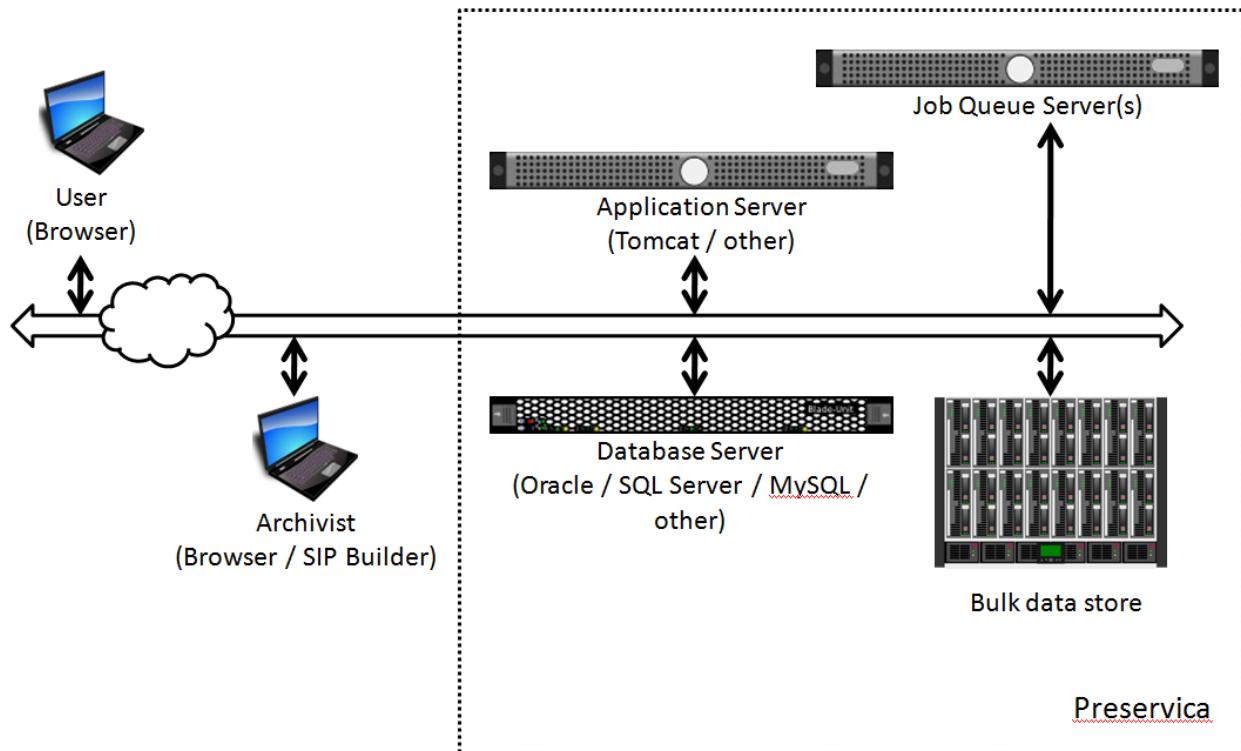
Tessella has developed Preservica, a Software-as-a-Service product, which has evolved from its Safety Deposit Box (SDB) application. Preservica is a multi-threaded application that can be installed across multiple server systems. Using a variety of hardware setups, Tessella has performed a number of studies to measure performance. These have shown that throughput is typically limited by the speed at which content can be supplied (which is often restricted by either network speeds or even the speed at which content can be read from storage). With the right environment ingest rates in excess of 20TB per day can be achieved routinely [14].

In addition, the human resources of an organisation are limited. Thus, to avoid this becoming a bottleneck, it is important that ingest, access and preservation workflows can be automated as much as possible. To mitigate this issue SDB utilises automated workflow steps wherever this is possible. This is also one of the unique features of the SDB information model: it allows for automated preservation decision making to occur “on the fly” based on human-entered policy information.

Figure 1 below demonstrates a typical SDB deployment diagram. Additional services such as User Administration (LDAP / Active Directory) and email (SMTP) servers may also be required to facilitate secure login and notifications of workflow successes and/or failures.

---

<sup>2</sup> <http://www.tessella.com>



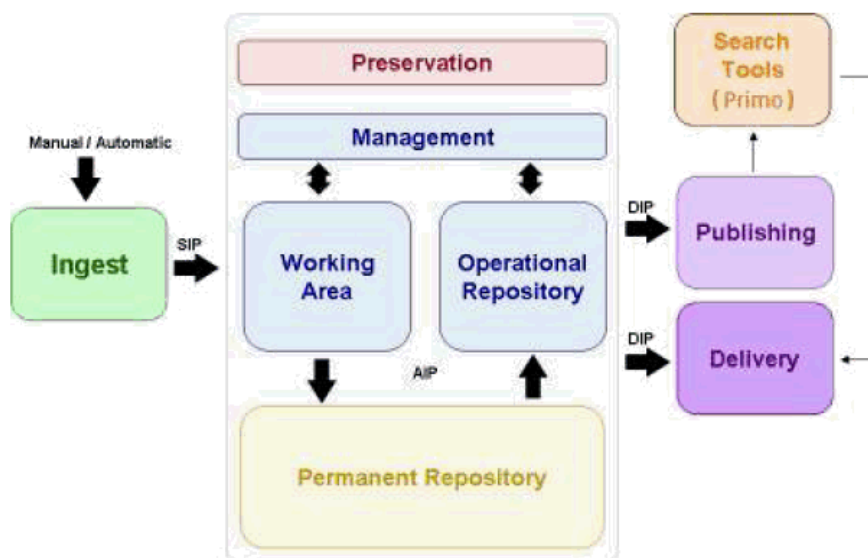
**Figure 1 – SDB / Preservica Physical Deployment Diagram**

The Preservica system allows an administrator to add in additional Application server nodes to facilitate higher rates of throughput. The Jobqueue Server is where the bulk of the Archive workflows are run and performs the “heavy lifting” of the system. The Database can be split in to shards to improve performance, and the bulk data store can be provided in multiple forms accessed in parallel.

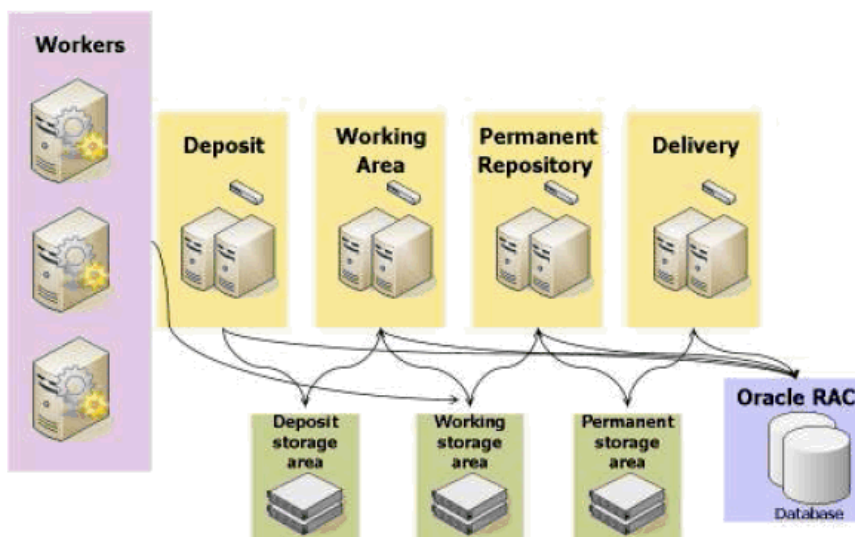
### 2.2.1.2 *Ex Libris Rosetta*<sup>3</sup>

Ex Libris Rosetta is a digital-object preservation solution that conforms to the ISO-recognized Open Archival Information System (OAIS) (see [1] [2] ) and supports international industry standards such as the Metadata Coding and Transmission Standard (METS), Preservation Metadata: Implementation Strategies (PREMIS), Dublin Core, and the Open Archives Initiative Protocol for Metadata Harvesting (OAI-PMH). The system is designed to support the acquisition, validation, ingest, storage, management, preservation, and dissemination of different types of digital objects and adheres to e-legal deposit requirements. Ex Libris Rosetta consists of a set of separate yet interactive modules that correspond with the life cycle of a digital object. The Ex Libris Rosetta system architecture is shown in the following diagram:

<sup>3</sup> <http://www.exlibrisgroup.com/category/RosettaOverview>



Ex Libris Rosetta provides a scalable infrastructure to address the ever-growing need to preserve and manage digital materials. This scalable solution can be implemented on a robust distributed architecture which allows the deposit module, working area, permanent repository, and database to be deployed on separate servers. Each module can be scaled up with additional of computers. In addition, Ex Libris Rosetta enables institutions to add dedicated servers, called *workers*, to perform specific tasks, such as virus and fixity checks. This flexibility allows an institution to start with a small hardware configuration and expand Rosetta to meet the needs of the institution's growing collection. Beyond the benefit of scalability, the architecture provides institutions with a system that offers redundancy and has no single point of failure. See the Scalability diagram below:



For a report on a scaling proof of concept see [17] .



## 2.2.2 Open Source

### 2.2.2.1 Duraspace<sup>4</sup>

Duraspace is an open source umbrella for several projects and services in the area of long term digital preservation that is used in more than 1000 repository instances in 90 countries. It includes five projects and each one of them has its own scalability considerations. The projects are:

- DSpace – institutional repository application
- Fedora – framework for building digital repositories
- Vivo – semantic web platform enabling research discovery
- DuraCloud – service for archiving content in the cloud
- DSpaceDirect – service for hosted DSpace repository solution

The DSpace application shows constant improvement in its scalability and performance. With so many deployments, it can study the bottlenecks and resolve them over time. DSpace is built on top of PostgreSQL database and there were some performance issues caused by bottlenecks in the way DSpace used this database. Cambridge institutional repository (<https://www.repository.cam.ac.uk/>) reported that they archived over 100,000 items and performance degradation started from couple of tens of thousands of items. They identified memory leaks, inefficient database design and exhausted database pools. This also hampered availability as upload of new items or even batch upload became very slow. DSpace also utilizes a search engine underneath and the Imperial College reported they have archived more than 122,000 items with no performance problems but the indexing process has an extremely non-linear response to archive size, and may become unusable as the archive expands. However, as database technology becomes more scalable e.g., Postgres-XC (eXtensible Cluster) and search engine technology becomes distributed e.g. Elastic Search, so DSpace can take advantage of these new technologies and improve its scalability even further.

The Fedora framework performance and scalability was studied extensively by Max Planck Society and FIZ Karlsruhe in the eSciDoc project in 2008<sup>5</sup>. They tested ingest of 14 million digital objects (patents database) and reported that Ingest times have remained stable for all objects. However, the ingest process took three weeks total, and there was a desire to improve that. The recent Fedora Futures is the next stage of Fedora and will address the top priority requirements expressed by the international community. The first requirement concerns improved performance and enhanced vertical and horizontal scalability. More specifically, Fedora Futures has requirements for mass object creation; ingesting large files; supporting Hierarchical Storage Management (HSM). Also, there is a requirement for multi-tenancy in the same fedora instance, without having to resort to using naming conventions and XACML policies to restrict access.

Vivo is a relatively new project in Duraspace, so there isn't much published information about its scalability. Vivo's reliance on in-memory caching of RDF data had put limits on its ultimate scalability. As a result, it was converted to use Jena's SPARQL database (SDB) subsystem. SDB significantly reduces the baseline memory footprint, allowing Vivo installations to scale well beyond what has previously been possible. However, it's not clear yet whether there are further scalability issues.

DuraCloud and DspaceDirect are services on the cloud. They provide flexibility and scalability offered instantly as a service. DuraCloud stores the data in Amazon S3 cloud and can offer an additional copy in Amazon Glacier, San Diego Supercomputer Cloud Store, and Rackspace Cloudfiles. DuraCloud thus can offer scalability similar to Amazon S3. DSpaceDirect includes the DSpace application on the cloud and thus its scalability is presumably similar to that of DSpace.

---

<sup>4</sup> <http://www.duraspace.org>

<sup>5</sup> <http://fedora.fiz-karlsruhe.de/docs/Wiki.jsp?page=Scalability>

### 2.2.3 EU Research Projects

Some of the current EU research projects on digital preservation that consider scalability in their objectives:

#### 2.2.3.1 *ENSURE*<sup>6</sup>

Enabling kNowledge Sustainability Usability and Recovery for Economic value (ENSURE) is an EU FP7 research project. Guaranteeing long term usability for spiraling amounts of data produced or controlled by organizations with commercial interests is quickly becoming a major problem. Guided by real world use cases in health care, finance, and clinical trials, ENSURE extends the state of the art in digital preservation through innovative solutions considering:

- Cost and Value: Evaluate the cost and benefit of different quality solutions, enabling a business to choose the most cost effective solution
- Preservation Lifecycle Management: Build on industry standard lifecycle management approaches to manage the preservation lifecycle, meet regulatory compliance, allow changes in the preservation approach to reflect environmental changes, address evolution of ontologies and manage the quality of digital objects over time
- Content-aware Long Term Data Protection: Provide data protection over long periods of time, addressing changes to personally identifiable information, new and evolving regulations, and manage user identities over the decades
- Utilize Emerging ICT: Evaluate the costs, risks and benefits and demonstrate how to use emerging, commonly available Information Technology to enable scalable solutions for digital preservation, in particular considering cloud storage and virtualization techniques.

The approach to scalability in ENSURE is in two different areas. The first is by considering scalability throughout the design of the system, build architecture that is scalable and can support large amount of complex data, taking into account constant interaction with the system from a distributed environment. The other area is looking at the scalable Information Technology solutions and the way to use them for digital preservation, such as using multiple cloud models for preservation and performing preservation activities by the storage solution to avoid moving data.

#### 2.2.3.2 *SCAPE*<sup>7</sup>

Scalable Preservation Environments (SCAPE) is an EU-funded project which is directed towards long term digital preservation of large-scale and heterogeneous collections of digital-objects. Its aims are to develop scalable services for preservation planning and preservation actions on an open source platform. These services will be based on a framework for automated, quality assured work-flows, which will be elaborated and tested during the project runtime. A policy-based preservation planning tool and an automated watch system will ensure a secure and targeted implementation of institutional preservation strategies.

SCAPE preservation components will be able to:

- Identify the need to act to preserve all or parts of a repository through characterization and trend analysis
- Define responses to those needs using formal descriptions of preservation policies and preservation plans
- Allow a high degree of automation, and scalable processing
- Monitor the quality of preservation processes.

Scalability is one of the scientific and technical objectives of the project. SCAPE addresses scalability in four dimensions:

---

<sup>6</sup> <http://ensure-fp7.eu>

<sup>7</sup> [www.scape-project.eu/](http://www.scape-project.eu/)



- Number of objects
- Size of objects
- Complexity of objects
- Heterogeneity of collections

The project is concerned with extending repository software functionality to enable storage, management, and manipulation of very large objects (e.g., multi-gigabyte video streams) and very large numbers of objects (hundreds of millions). Improving the ability of a number of existing preservation tools to manage a variety of container objects and to recognize diverse object formats is another SCAPE aim.

### 2.2.3.3 SCIDIP-ES<sup>8</sup>

The EU funded Science Data Infrastructure for Preservation - Earth Science (SCIDIP-ES) FP7 project is developing a data preservation infrastructure for scientific data using Earth science data as its testbed.

The aim of the initiative is to deliver generic services for science data preservation as part of the data infrastructure for e-science and to build on the experience of the ESA Earth Observation Long Term Data Preservation (LTDP) program to favour the set-up of a European Framework for the long term preservation of Earth Science (ES) data through the definition of common preservation policies, the harmonization of metadata and semantics and the deployment of the generic infrastructure services in the ES domain. The generic preservation services will address science data persistent storage, access and management and will build on the already proven research prototype services from the FP5 CASPAR project. They will be evaluated and tuned in depth, using Earth Science as pathfinder, and broadly but less deeply across other disciplines. The generic services will allow all kinds of data to be usable by researchers from many different domains and will enable the cost for long-term usability across disciplines to be shared supporting the creation of strong business cases for the long-term support of that data.

SCIDIP-ES aims at:

- Delivering generic sustained services and toolkits for long-term preservation and usability as part of the data infrastructure for e-Science
- Harmonizing data preservation policies, approaches and tools in the Earth Science Domain. Demonstrating the services through advanced test-beds.

Scalability of the toolkits and services used is part of the project's objective. A report was delivered in May 2012 (see [7] ) that analyzes the scalability aspects of the tools and services and recommends a plan to enhance scalability.

## 2.3 CLOUD STORAGE

### 2.3.1 What is Cloud Storage?

One of the key advances in storage solutions today is the development of *Cloud Storage* technologies, which provide data and storage support in *The Cloud*. Cloud storage refers to a virtualized entity of data that is available and accessed *online* and hosted on a variety of multiple virtual servers, mostly hosted by third parties, rather than being hosted on dedicated servers. The data is distributed and is typically hosted in geographically dispersed locations. It facilitates the storage utility as a service, where storage is offered in a pay-per-use model.

From a technological perspective, cloud storage provides new levels of scalability and elasticity and availability<sup>9</sup>, and enables simple access to data from any location and any device. Virtualization is key – the user is not aware of where the data is being stored and how it's being managed; this greatly

---

<sup>8</sup> <http://www.scidip-es.eu/scidip-es>

<sup>9</sup> As of April 2013, Amazon's S3 stores 2 trillion objects, the largest cloud storage provider today.

simplifies the use of storage and reduces its total cost of ownership (TCO). It is an infrastructure as a service (IaaS) that can be used as primary storage by Web 2.0 developers and as a low-cost alternative storage tier for enterprises and small and midsize businesses (SMBs). Storage cloud infrastructure is increasingly positioned as a minimally functional low-cost tier of storage to businesses and enterprises, usually, but not always, positioning cloud storage for latency-tolerant applications such as backup and archiving, thus making it attractive for digital preservation repositories.

Probably the biggest attraction of cloud storage is the price. The pricing is composed of a monthly fee per-GB, with additional fees every time data is moved into or out of the service provider's storage environment. This pricing scheme makes it well suited to applications such as archiving and long term preservation, which have relatively low access activity. Furthermore, the pay-as-you-go model has turned the cost from capital expenditure into operations expenditure, thus making it economically attractive to many institutions.

At its early stages, cloud storage services provided very basic storage functionality – primarily storing and retrieving objects. Today, richer functionalities are being offered, including, for example, varying levels of reliability, deletion, retention, and logging.

A highly relevant development in the area of Cloud technologies is the emergence of a new open source project called *OpenStack*<sup>10</sup>. *OpenStack* is an Apache 2.0 licensed software for building private and public clouds. It is a global collaboration of developers providing cloud computing (*IaaS*) *open source* platform for compute, network and storage. The *OpenStack Swift* project is a highly-available, distributed object store based on commodity hardware. OpenStack is gaining a lot of momentum and its Swift solution may become a real relevant option for building cloud-scale storage for archiving and preservation.

In summary, the value proposition of cloud storage technologies lies in their scalability in the number of objects and in the aggregated quantity of the data, the low initial and fine-grained incremental costs, the availability anywhere and the ability to share data across domains and geographies. However, it was observed early on that simply "throwing" data on a cloud is not a good solution for digital preservation repositories [12] [13] and more advanced management and reliability mechanisms are needed. It is yet to be seen whether storage cloud technologies will evolve in a manner that will make them suitable for digital preservation in the future.

One of the ways to evolve cloud storage in a manner that will serve the preservation community as a whole is through standards and specifications. Specifically, the CDMI (Cloud Data Management Interface) [11] standard, developed by the Storage Networking Industry Association (SNIA), is a functional interface that applications use to store and access data elements in the Cloud. As part of this interface the client is able to discover the capabilities of the cloud storage offering and use this interface to manage containers and objects, including setting of metadata on these data elements. CDMI can also be used to manage accounts, security, billing, etc, even for storage that is accessible by other protocols. The CDMI specification is a SNIA Architecture standard.

### 2.3.2 Scalability in Cloud Storage

As noted earlier, scalability is a key *proven* characteristic of the Cloud technology. Indeed, Cloud addresses the growing demand in capacity and availability in an unprecedented way. If we inspect the capacity of data stored in public clouds today, we observe that cloud providers succeed in growing persistently in terms of total capacity, number of objects stored and number of users.

We provide below some indications on the scalability of existing public clouds and OpenStack.

#### 2.3.2.1 Amazon S3<sup>11</sup>

Some published facts about Amazon Web Services (AWS) S3 scalability:

- A report from June 12, 2012<sup>12</sup> :

---

<sup>10</sup> <http://www.openstack.org/>

<sup>11</sup> <http://aws.amazon.com/s3/faqs>

- "The number of objects stored in Amazon S3 reached one trillion (1,000,000,000,000 or  $10^{12}$ ). Lately, we've seen the object count grow by up to 3.5 billion objects in a single day (that's over 40,000 new objects per second)"
- Since AWS introduced S3's object expiration feature, its customers have used it to delete over 125 billion objects. Still, the overall object count has continued to grow very rapidly, although deletion has been made very easy
- The total volume of data and number of objects you can store are unlimited
- Individual Amazon S3 objects can range in size from 1 byte to 5 terabytes
- The largest object that can be uploaded in a single operation is 5 gigabytes. For objects larger than 100 megabytes, customers should consider using the Multipart Upload capability. SLA does not measure bandwidth – only "Uptime percentage" and "Error rate"
- **Capacity limitations:** 5 TB on single object. No limitation on container
- **Bandwidth limitations:** No reference to Bandwidth limitations
- **Availability (time to access, access pattern):** SLA for Uptime of 99.9%
- **Distribution (data location, access locations)** Amazon S3 offers storage in the US Standard, US West (Oregon), US West (Northern California), EU (Ireland), Asia Pacific (Singapore), Asia Pacific (Tokyo), Asia Pacific (Sydney), South America (Sao Paulo), and AWS GovCloud (US) Regions. You specify a Region when you create your Amazon S3 bucket. Within that Region, your objects are redundantly stored on multiple devices across multiple facilities
- **Multi tenancy:** Open to all customers. Claim that Amazon uses S3 for internal use as well. Limited ACL support to accommodate
- **Long-term aspects:** *Glacier service*<sup>13</sup> is targeted for archival with long access times and greatly reduced cost.

#### 2.3.2.2 CloudFiles Rackspace<sup>14</sup>

Some published facts about Rackspace cloud service:

- **Capacity limitations:** Rackspace provides "unlimited online storage" and CDN services. Uploading content to Cloud Files through the Control Panel is limited to files below 5GB; transfer of files larger than 5GB is accommodated by enabling segmentation of files into multiple file segments
- **Container Size:** For best performance, the container size is recommended to have no more than 500,000 objects
- **Bandwidth limitations:** There is no reference to bandwidth limitations. The service can be combined with Akamai CDN services for improved accessibility on reads. "Bulk upload" can be done by sending in a hard disk by mail
- **Availability SLA:** "... guarantee that Cloud Block Storage service will be available 99.9% of the time in a given billing cycle"
- **Data Distribution:** Multiple copies in zones in the same Data Centre. Several sites in US, and one in London
- Supports **Multi tenancy**
- No special long term services / aspects.

#### 2.3.2.3 Google Cloud Storage<sup>15</sup>

Some published facts about Google Cloud service:

<sup>12</sup> <http://aws.typepad.com/aws/2012/06/amazon-s3-the-first-trillion-objects.html>

<sup>13</sup> <http://aws.amazon.com/glacier/>

<sup>14</sup> <http://www.rackspace.com/>

<sup>15</sup> <https://cloud.google.com/products/cloud-storage>

- **Capacity limitations:** Unlimited. Pricing becomes cheaper with higher usage. Changes according to geography.
- **Availability SLA:** with 99.9% availability as base.
- **Distribution:** data is replicated in multiple data centres that are geographically distributed for greater availability. Data is also replicated within each data centre. Currently, data is stored in U.S. and Europe data centres only
- Supports **Multi tenancy**.

#### 2.3.2.4 iCloud<sup>16</sup>

iCloud seems highly scalable based on the amounts of data it stores.

A recent study by Strategy Analytics noted that Apple's iCloud is now the largest supplier of online media services for American consumers, equaling the usage of Dropbox and Google combined."

- The company's data centers also handle massive terabytes of iCloud documents and Photo Stream pictures, iOS device backups, emails, calendars and contact synchronisations, as well as supporting iPhone location lookups, iBooks, Maps, Siri and software updates for Mac OS, iOS and hundreds of thousands of App Store titles
- In April 2013, Apple reported coordinating the billionth download of free iTunes U educational content
- As of April 2013, the service has 300 million users
- No SLA guarantees for **availability**
- **Distribution:** services are hosted on several large Data Centers in the US. Expanding its infrastructure and building new data centres
- Supports **Multi tenancy**.

#### 2.3.2.5 DropBox<sup>17</sup>

Drobox works via Amazon S3, and as such inherits the scalability properties of AWS S3. On 12 November 2012, Dropbox announced it had reached 100 million users. It supports multi tenancy.

#### 2.3.2.6 OpenStack Swift

OpenStack Swift is scalable software; hence it can be configured and deployed at various levels of scalability.

- **Capacity limitations:** limited to single location due to communication speed issues. Within a single cluster have not found limits. Supposed to scale-out horizontally. "The San Diego Supercomputer Center (SDSC) has set up the largest academic cloud storage using OpenStack Swift with 5.5PB of raw storage (1.8PB-2.2PB usable given two-way replication"
- This service is designed for high availability (based on multiple replication)
- **Distribution:** all data is stored within a single cluster. In July 2013, SwiftStack announced global clusters support for Swift 1.9.0<sup>18</sup>
- Supports **Multi tenancy**

### 2.3.3 Cloud Delivery Models

There is a tendency to speak of "cloud" as if all clouds are created equal. They are not. When we say something is a cloud, all we mean is that we "enable convenient, on-demand network access to a shared pool of configurable computing resources that can be rapidly provisioned and released with minimal management effort"[8] . These resources can be delivered in a range of ways, each of which has various advantages and disadvantages.

---

<sup>16</sup> <https://www.icloud.com/>

<sup>17</sup> <https://www.dropbox.com/>

<sup>18</sup> <http://swiftstack.com/blog/2013/07/02/swift-1-9-0-release/>

While there are various sub-variants, there are three main delivery models: public, community and private.

A *public cloud* is one in which anyone or any organization can gain access, and the resources are shared among all users of the cloud. The best known public cloud for infrastructure services is Amazon Web Services; however, there are many other providers such as Rackspace, IBM and Microsoft. Public clouds are accessed over the Internet. Public clouds provide the most elasticity and the most flexibility in rapid provisioning. In a public cloud, a user has no idea and no control over whom else is using these resources. Furthermore, public clouds have SLAs and management processes defined by the cloud operator and as a consumer of the services, there is little to no room for negotiation and often these processes may be opaque. Because of this, public clouds may be best suited for use in preservation either if multiple clouds are used [9] or as a backup archive [10] .

*Private clouds*, at the other extreme, are clouds operated for the sole use of one organization; this organization may have departments which share the private cloud. A private cloud is typically accessed by an intranet. Since a private cloud is being run for only a specific organization, that organization can own and control the management processes and specify the desired SLAs. This said, private clouds are not for everyone. Elasticity and the ability to rapidly respond are a function of the size of the cloud and the number of users. Thus, private clouds are limited in their elasticity and rapidity of provisioning except when deployed by large organizations, e.g., national governments.

*Community clouds* are run for a consortium of cooperating, like-minded organizations, e.g., a collection of universities or libraries. Community clouds can have the benefits of both public and private clouds. By aggregating a set of smaller organizations with similar objectives, the community cloud is able to be large enough to have elasticity and rapid provisioning, while still implementing processes specified by the community and in a manner that is transparent to the community. The drawback of community clouds is the need for all of the participating organizations to trust one another; thus, in general, community clouds will be more appropriate to use for preservation of cultural heritage data than for commercial data.

As should be clear, there is no one "right" answer in terms of the appropriate cloud delivery model. For a secondary copy, public clouds may make the most sense, for very large organizations, a private cloud could be optimal, and for like-minded organizations, a community cloud could be most appropriate.



### 3 SCALABILITY ANALYSIS

#### 3.1 GENERAL REMARKS

In this chapter we analyze scalability. We start by listing the parameters that affects scalability of digital preservation. Then we present the results of a questionnaire that included questions on scalability and was distributed to selected organizations maintaining digital repositories. The questionnaire includes questions related to the repositories nature, storage solutions, and scalability. In this report we concentrate on the scalability aspects of the repositories and owning organizations. Analysis of the results follows.

The results presented are based on 18 answered surveys in Italy and Germany by European Space Agency (ESA), Consorzio Interuniversitario Nazionale per l'Informatica (CINI), Deutsche National Bibliothek (DNB), plus additional information from a different survey held in the Netherlands by National Library of the Netherlands (KB) in the 1st half of 2012. In addition, we used data from a PARSE.Insight survey held in 2009 (see [15]) with answers from thousands of respondents.

Some of the questionnaires were only partially answered. In the detailed analysis below, we explicitly refer to the number of missing responses as it is meaningful to the generalization of the results.

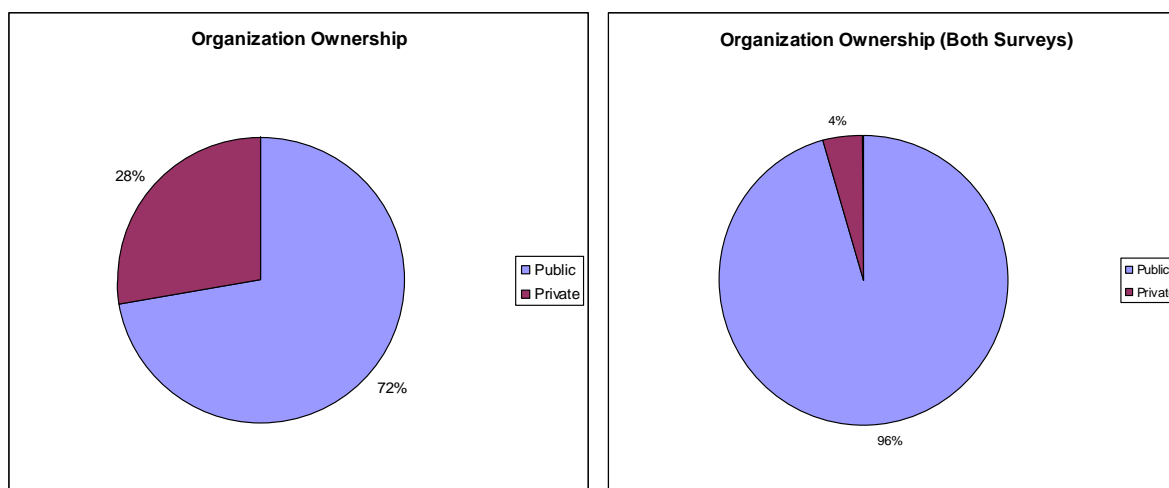
#### 3.2 SCALABILITY PARAMETERS

We have identified a list of parameters that are important to scalability:

1. **Variety of data types** – the more data types there are to preserve, the more scalable a preservation system is required, as each data type adds unique requirements to the system in order to support its preservation, such as understanding the metadata, indexing, representation information, packaging to AIP, preservation actions, and retention policy.
2. **Size of the data** – the total capability of the data preserved sets the scalability requirements for the preservation storage. It also affects the amount of processing required for the preservation actions, the network capacity to support data transfers, and the size of the search indexes.
3. **Number of digital objects and size of each object** – higher number of objects preserved requires more scalable mechanisms as the system needs to be able to maintain and access each object. Larger varieties in object sizes (e.g. video objects vs. small documents) require mechanisms that support all those sizes
4. **Number of versions of the same object** – multiple versions require mechanisms to support maintaining the connection between versions during the data life cycle
5. **Connections between different objects** – the ability to connect objects, e.g. connection between a publication and the underlying data it uses, enhance the usability of the system, however maintaining a large number of connections require scalable mechanisms in the system
6. **Amount of metadata and connections between metadata** – more metadata means more data is preserved. In addition maintaining the connections between the data and metadata and amongst metadata requires more scalability of the system
7. **Retention Period** – longer retention periods require support for more accumulated data and more preservation actions
8. **Availability** – are objects are expected to be available at any time from anywhere? Higher availability requires more resources
9. **Distribution** – how geographically dispersed the system is affects scalability. Supporting distributed systems is more complex
10. **Searchability** – maintaining the ability to better search for the preserved data requires better search indexes, and better response times
11. **Degree of sharing** – supporting multiple curators and multiple users in the system adds concurrency requirements to the system. Multi tenancy adds many requirements on security, capacity, support for multiple policies, SLA and more..

### 3.3 RESPONDING ORGANIZATIONS

The majority of the organizations responding to the survey are public organizations – 13 out of 18. The extended data from the survey by KB is from 87 more public organizations. The first chart below shows the percentage of public and private organizations responding to the survey carried out within this project. The second chart shows the percentage of public and private organizations when including the KB survey results with the project survey.

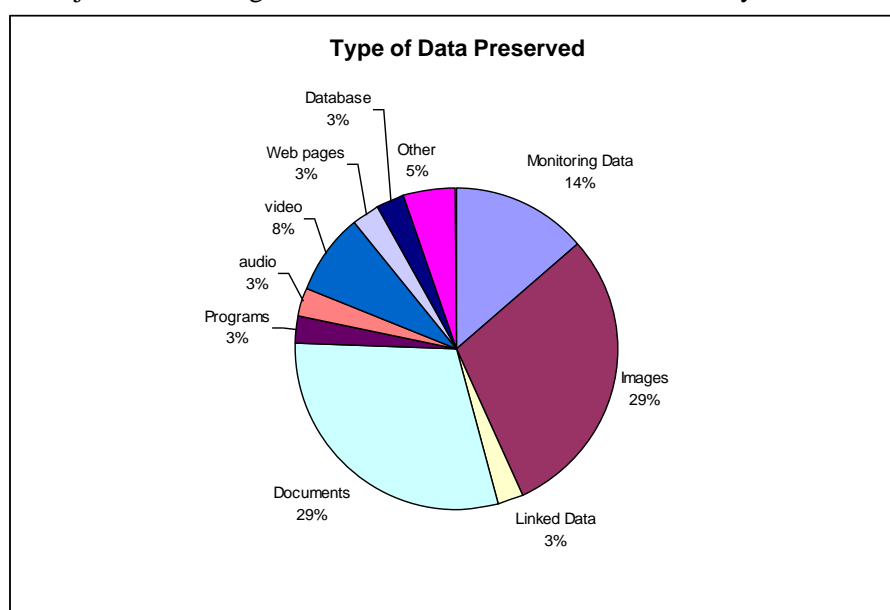


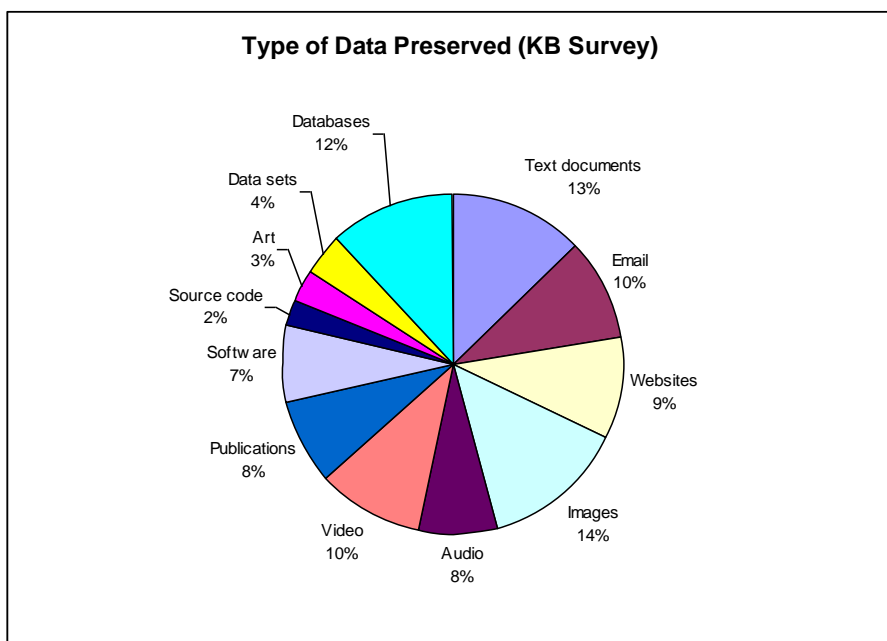
### 3.4 RESULTS AND ANALYSIS

We present the survey results for the scalability parameters identified in section 3.2 and discuss their implications.

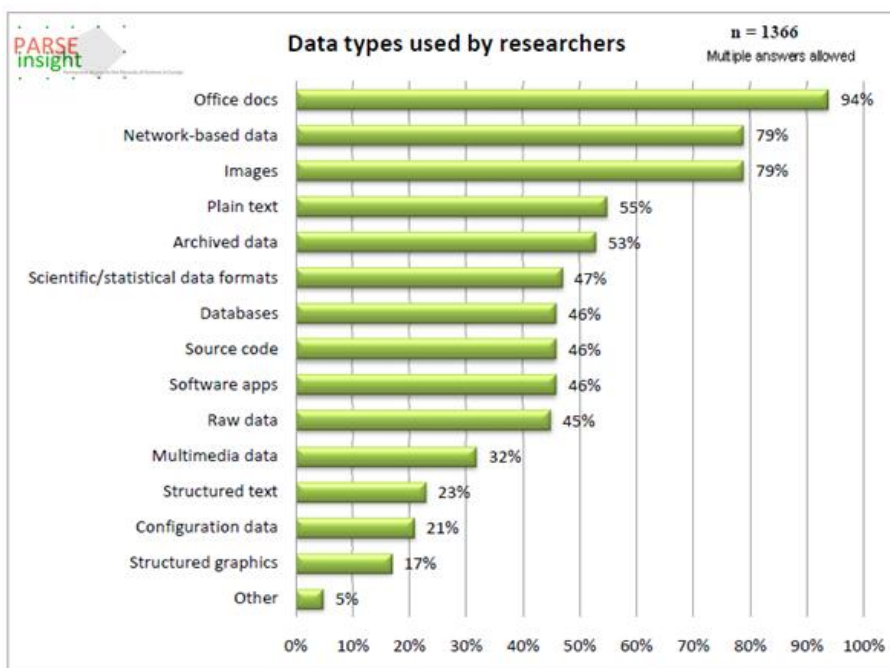
#### 3.4.1 Data Types Preserved

The figures below show the different kinds of data preserved from the APARSEN storage and scalability survey and from the survey by KB respectively. Each category in the diagrams aggregates multiple formats, so the overall variety is very large. The range of sizes is also enormous, e.g., monitoring data objects can be huge, while some text documents can be very small.

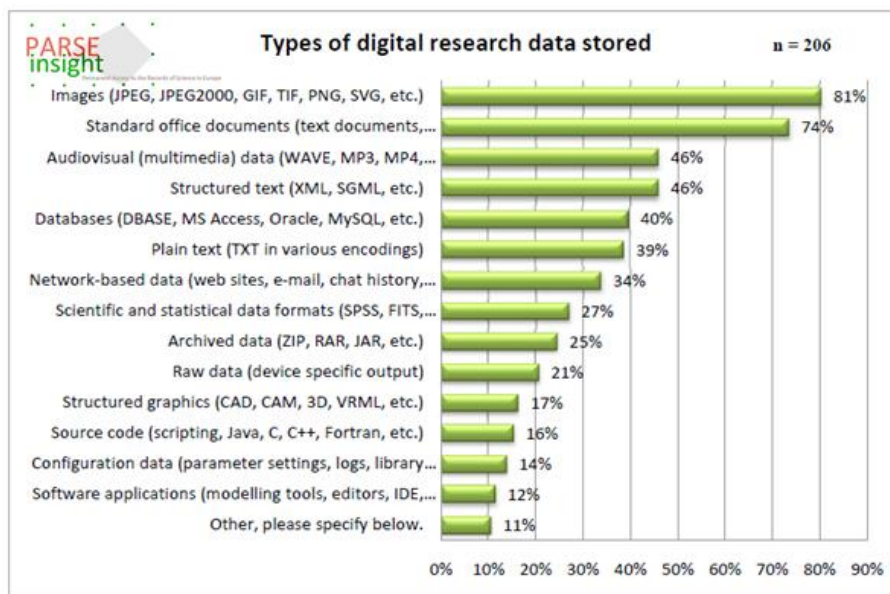




The survey results from PARSE.Insight are consistent with the above showing large range of data types and sizes:

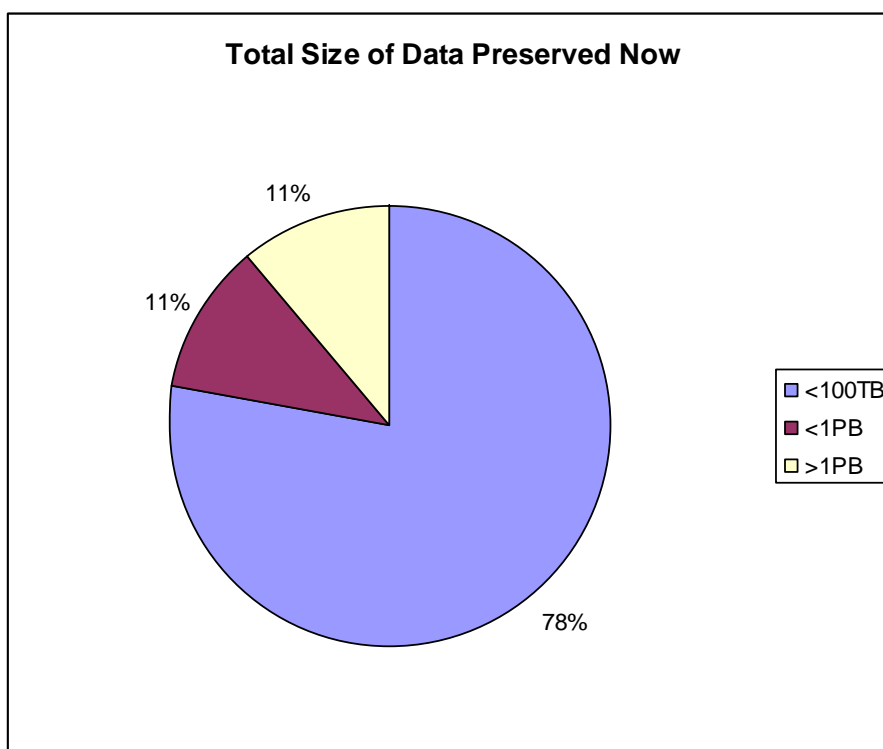




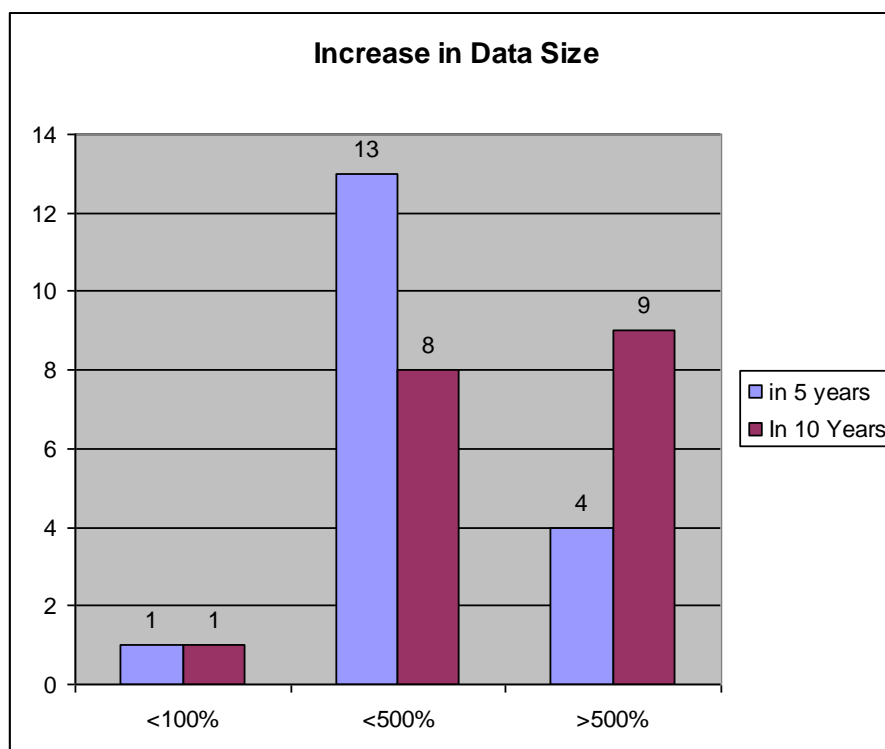


### 3.4.2 Size of the Data

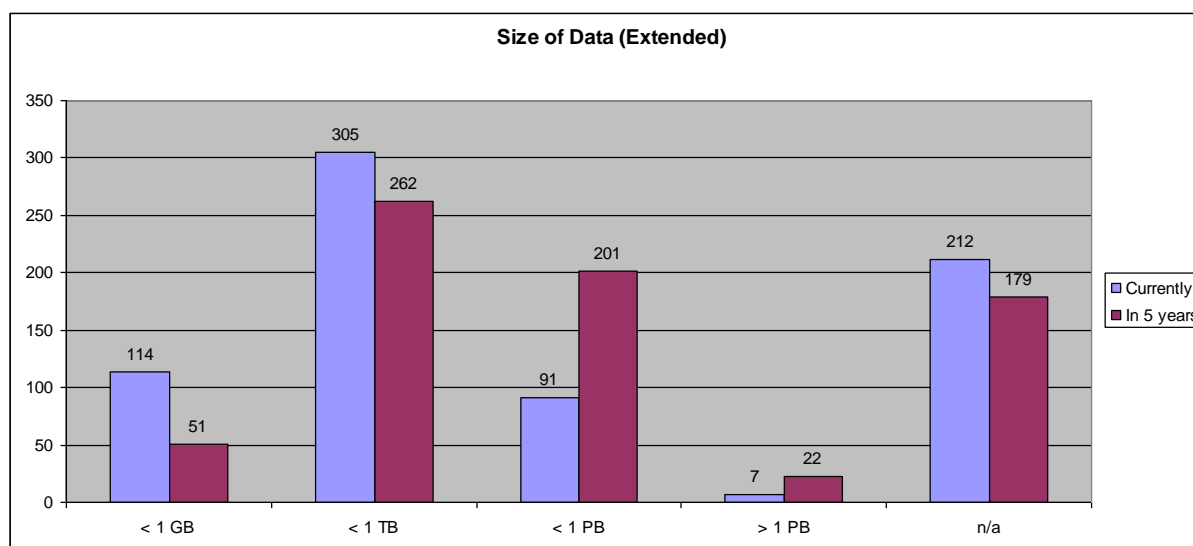
The survey results show most organizations expect significant growth to the current size of their repositories. The following charts give the current size of the data and the expected growth in the next 5 and 10 years.



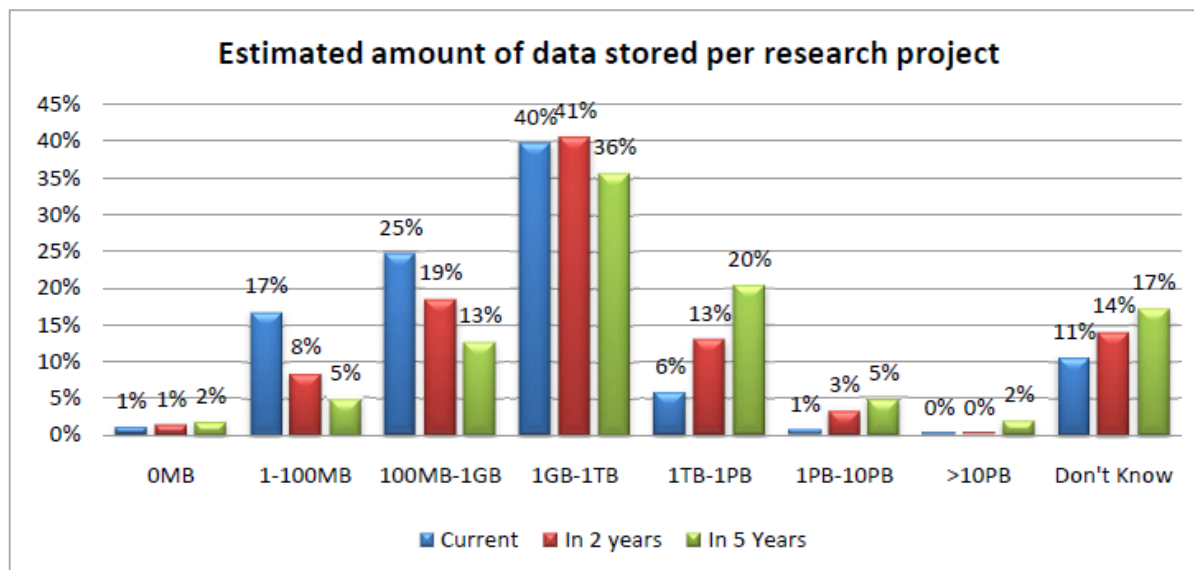
The chart below shows that most organisations (72%) expect data sizes to increase 5 fold in 5 years' time with almost half of organisations expecting data sizes to increase by more than 5 fold in 10 years' time.



The survey by KB shows similar trends with the majority of data sizes currently below 1TB and although in 5 years' time the majority of data sizes are still below 1 TB it is worth noting that data sizes below 1GB halve in 5 years' time whilst those above 1TB double in 5 years' time from current levels:



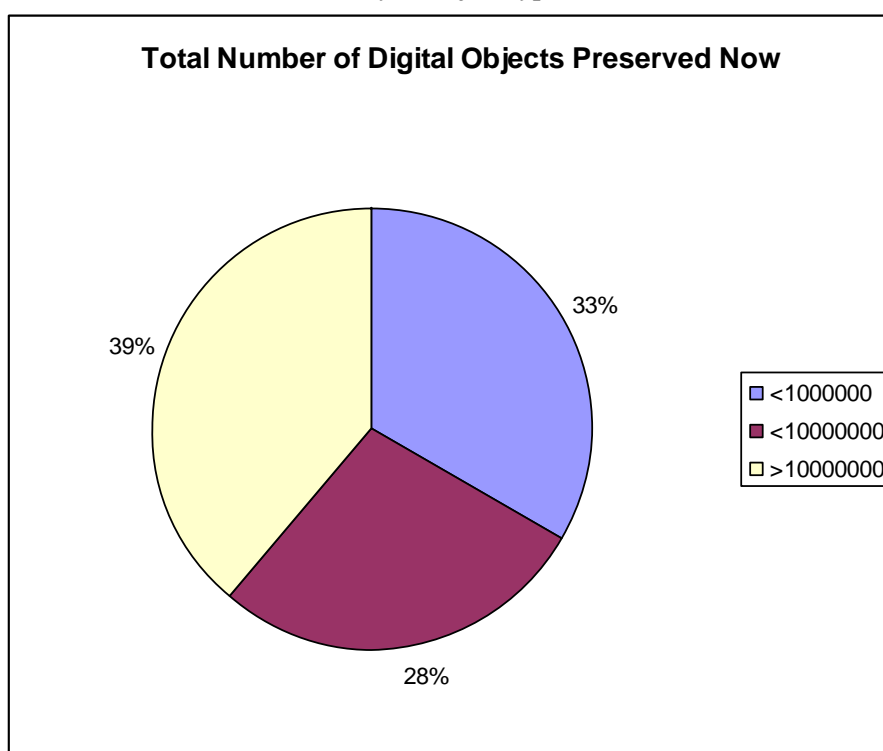
The PARSE.Insight data supports the trend of growth in the amount of data, showing increase in along the years. The percentage of projects (out of 1296 projects) with data of more than 1 Terabyte of data was forecasted to increase from 7% to 27% over 5 years..



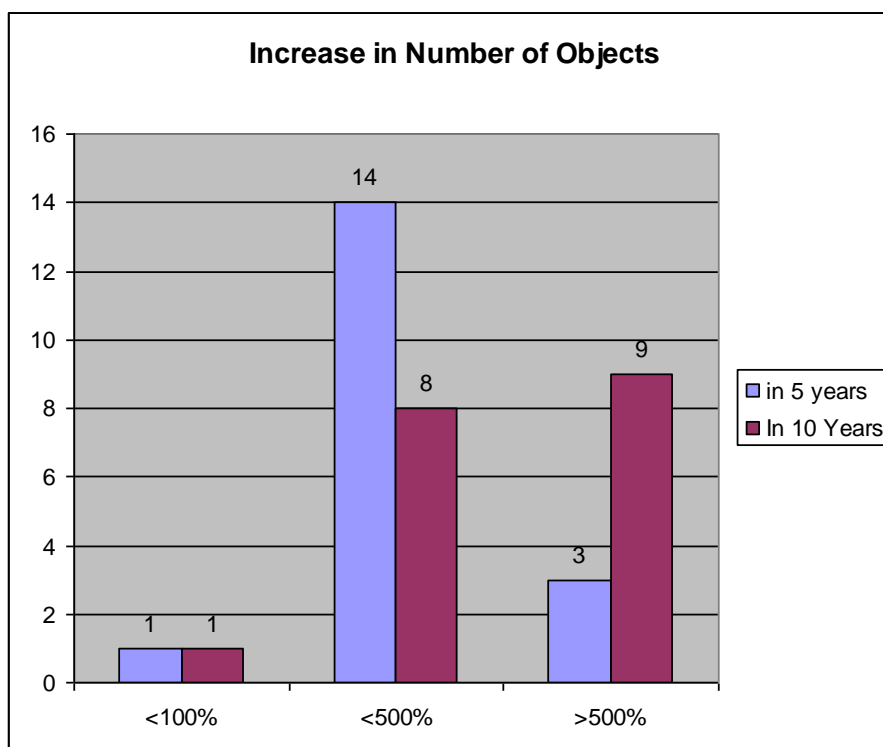
### 3.4.3 Number of Digital Objects and Their Sizes

The survey results show that most organizations expect significant growth to the number of digital objects they preserve. The following charts give the current number of objects and the expected growth in the next 5 and 10 years.

The variations of size follows from the variety of object types as discussed in section 3.4.1

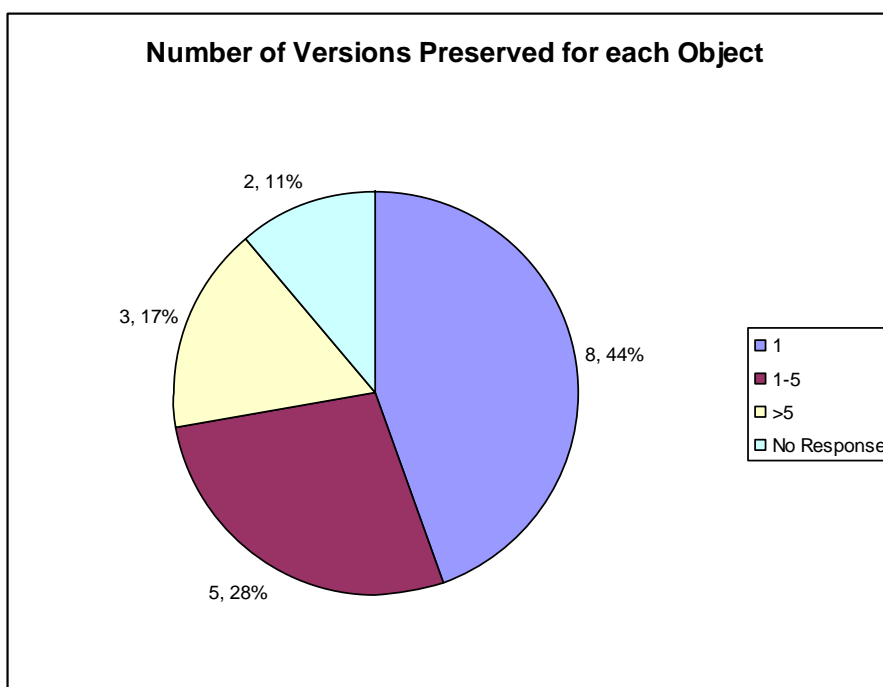


The chart below shows that the number of objects is expected to increase over both 5 and 10 years' time. The vast majority of the organizations (77%) expect growth of up to 5 fold and 16% expect more than 5 fold growth in 5 years. When asked about 10 years time, 50% are expecting more than 5 fold growth.

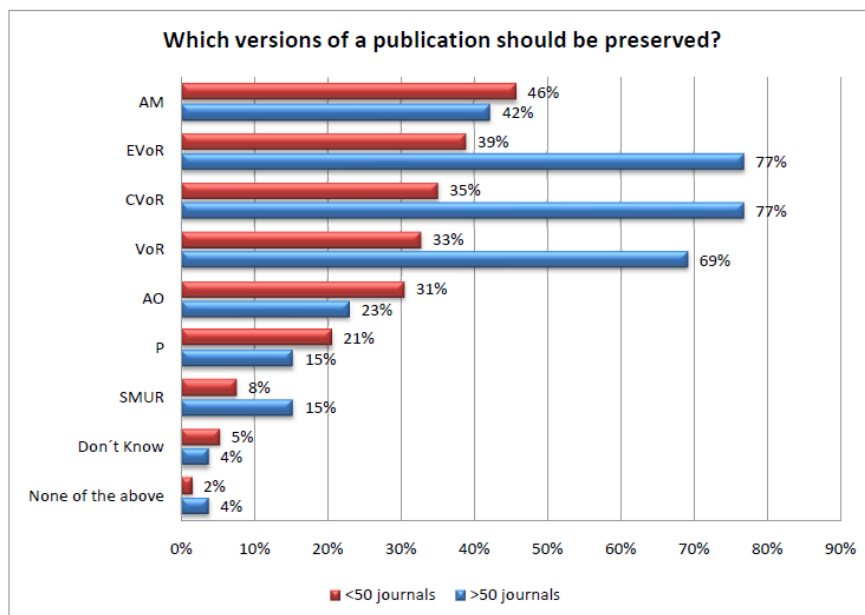


### 3.4.4 Versions

The survey results show that nearly half of the repositories keep only one version, however, 17% keep a large number of versions i.e. more than five versions of each object.

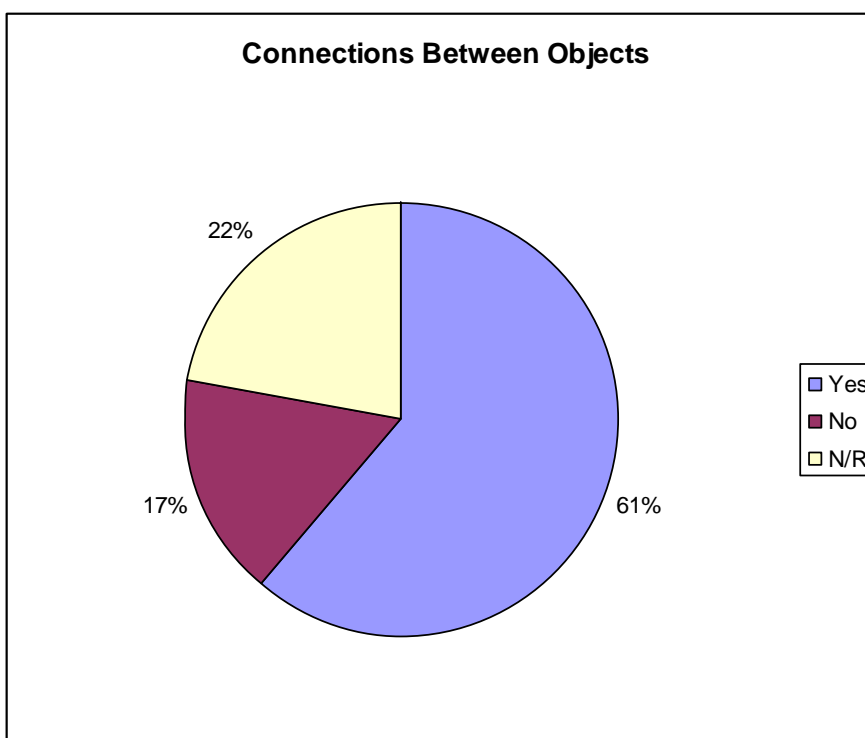


The PARSE.Insight survey results from publishers also show the need for multiple versions of publications:



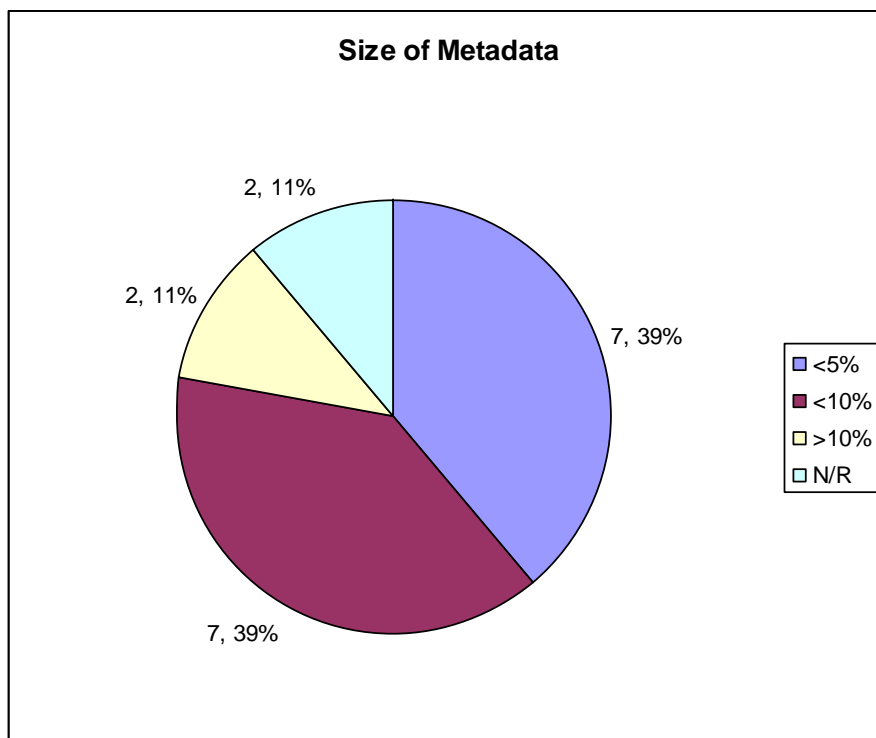
### 3.4.5 Connections between Objects

In most organizations there are connections between objects, increasing the complexity of the repository.



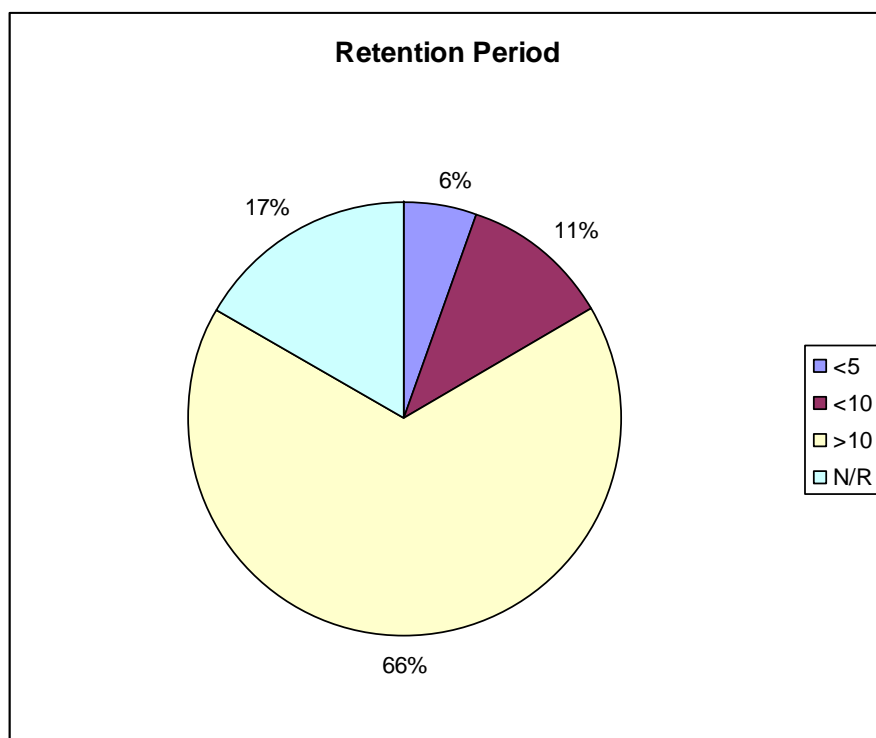
### 3.4.6 Size of the Metadata and Connections between Metadata

The survey results show that in most organizations the size of the metadata is not very high, only 11% indicate size of more than 10%. Connections between metadata were not surveyed, but we can infer from the connections between objects that such connections are likely to exist.

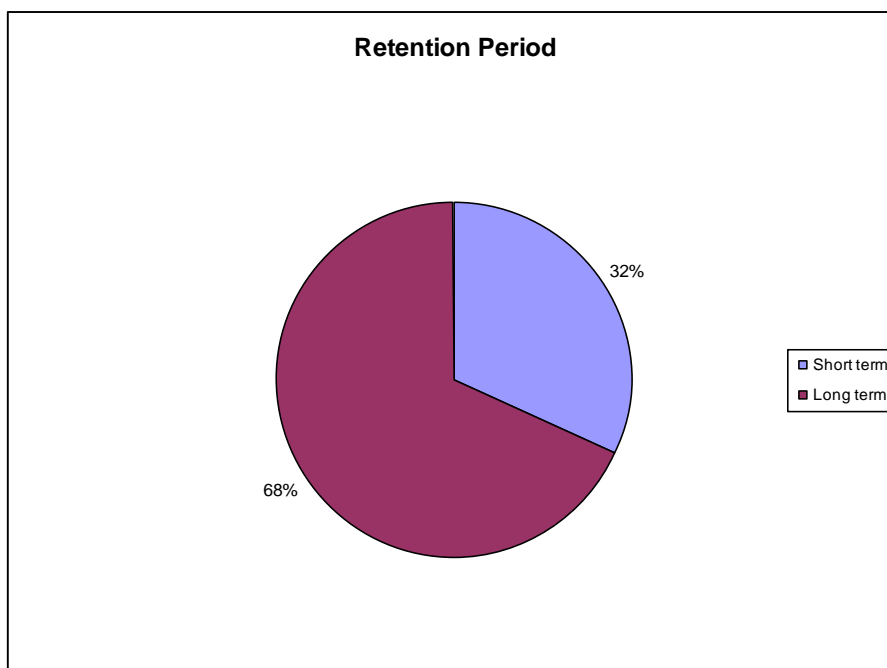


### 3.4.7 Retention Period

The survey results show that the vast majority of organizations (66%) preserve the data for more than 10 years.

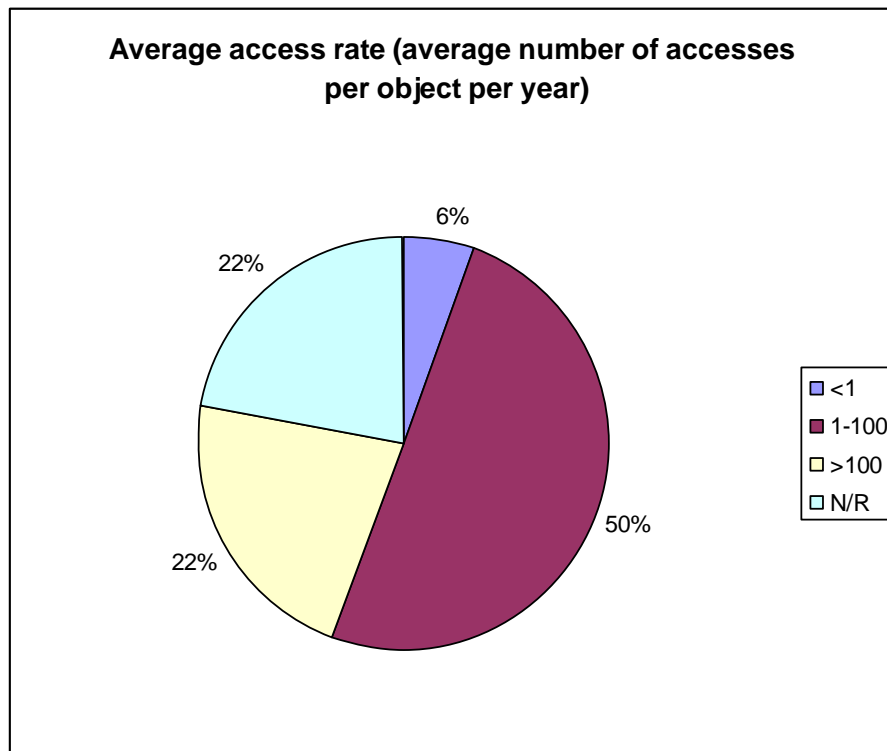


The survey by KB shows similar results, most data have a long term retention period.

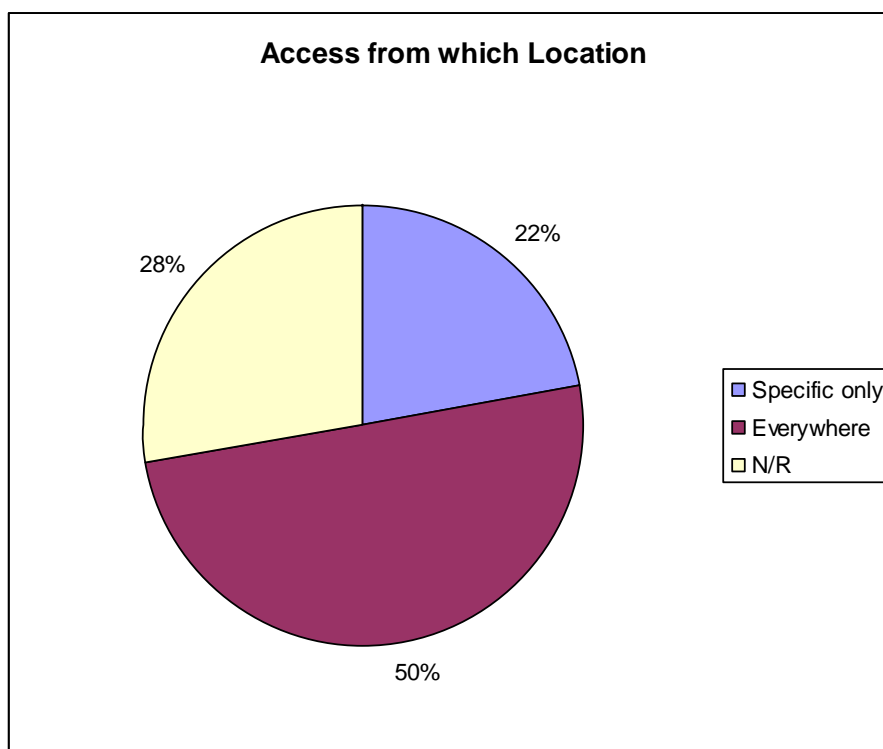


### 3.4.8 Availability (Access Rate, Access Location and Access time)

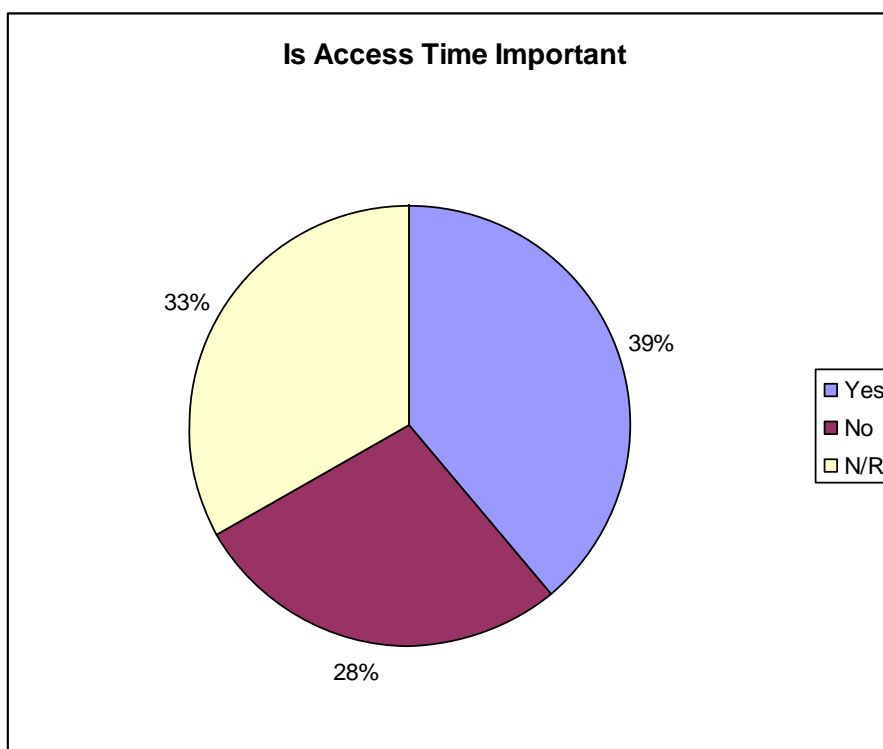
The survey results indicate that only a small fraction of organizations (6%) have very low access rates (less than once a year per object), nearly a quarter of the organizations (22%) indicated high access rates (accessing each object more than 100 times a year).



Half of the organizations indicated they require access to the data from anywhere.



Access time is important to 39% of the organizations whilst 33% responded that access time was not important.



### 3.4.9 Distribution

The survey did not ask about distributed repositories, however from some of the verbal answers, larger organizations often store data in more than one location. In the KB survey there was references to keeping data in a different country and 10% of the responding organizations responded they keep part



of their data abroad. Thus we can infer that at least part of the organizations have distributed repositories.

### 3.4.10 Searchability

The survey did not ask about searchability explicitly, however, several organizations stated that the data needed to be located based on varying access criteria and that response time was important, hence the importance of searchability is implied.

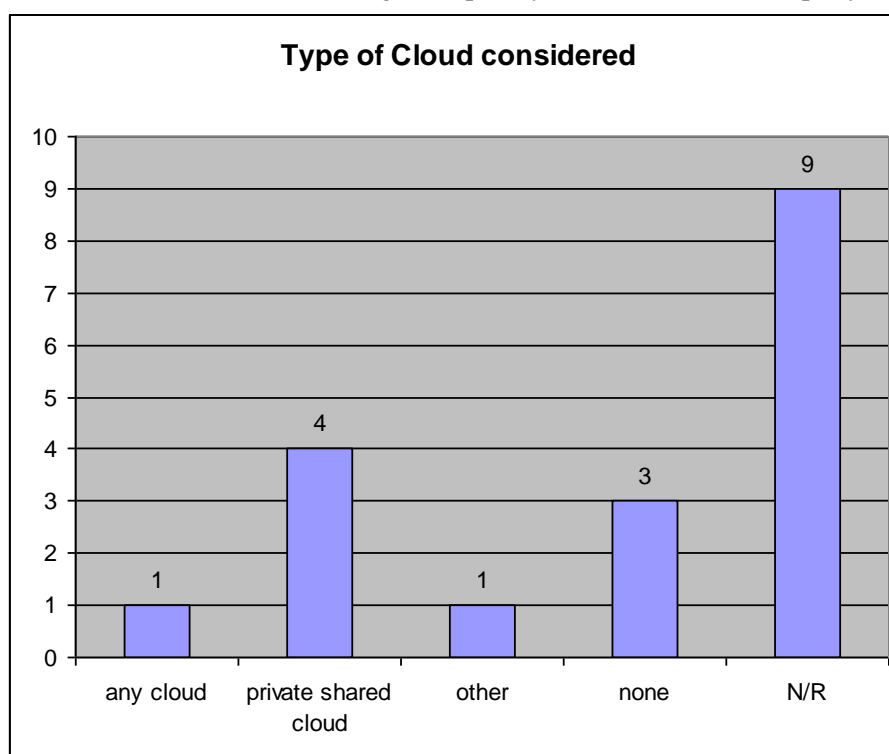
### 3.4.11 Degree of sharing

As all surveyed organizations except one use in-house storage, there is no multi-tenancy amongst those surveyed. One organization uses outsourced storage managed by an external company. This indicates that the degree of sharing is currently very low amongst the surveyed organizations.

The likelihood of moving was questioned, but only few answered. The results are summarized in the diagram below. From the 9 organizations that responded, 3 were not ready to consider the use of a cloud with multi-tenancy, 4 were ready to consider community cloud (private shared cloud), one was ready to consider any cloud, including public, and one did not identify a specific cloud delivery model.

When asked for reasons to consider cloud, the majority of those that answered mentioned scalability amongst other criteria such as reliability and availability. Amongst fear factors we can find concerns such as security, lack of control over the data, and vendor lock-in. Security concerns are rated high.

In the Survey by KB, all but 2 organizations, manage digital information, either for preservation, access or institutional management. Out of 87 respondents 33 (38%) manage their own storage solutions, whereas 39 (45%) partly manage their own solutions and have storage partly outsourced. Only 10%, that is 9 out of 87 have their storage completely outsourced to a third party solution.



### 3.4.12 Overall Results Analysis

Analyzing the survey results, we first note that the results we have are mostly from public organizations. The reason is that corporations are not keen to share such information and often keep it

confidential. Keeping this limitation in mind, the results we have here are valuable to identify gaps between the surveyed scalability factors and the required scalability by the same organizations.

Most of the repositories surveyed are not large, yet many of them expect considerable growth in the coming years. 50% expect growth of more than five fold in their size in 10 years. The expected increase in the number of preserved objects is similar, 50% expect growth of more than five fold in 10 years.

Most of the other parameters indicate that the complexity of the preservation data is high. There are a large variety of data types as indicated by both surveys. 50% of the repositories keep multiple versions, with some keeping a large number of versions. More than half (61%) indicated connections between objects. Combined with the result that about two thirds of the organizations require long term preservation of over 10 years, current preservation repositories are expecting some serious challenges ahead.

In addition, there other factors that present demands from the preservation system. Most organizations (72%) expect accessing the data multiple times with 22% over 100 accesses per object per year. Combined with the fact that half of the organizations want the access from anywhere in the organization as well as the importance of response times (39%), attention should be given to meeting these requirements when the size of the repository is due to grow substantially.

The vast majority of the surveyed organization use home grown solutions. Taking into account the complexity of the requirements combined with the expected growth in data size, there is a gap here. It will be hard to evolve those solutions to be more scalable when each organization has to cope with the development of the system and obtaining the storage by itself. As there are multiple dimensions to supporting a scalable system, creating the next level of scalable systems cannot be achieved by point improvements to non scalable systems. Design for scalability is required and it touches many aspects of the system.

Another gap relates to the scalability of the storage systems used. Most of the surveyed organizations use and maintain their own storage. With the growth in the size of data, this can become challenging, both with the need to obtain the storage and in maintaining it. Data from the PARSE.Insight survey support this observation, were 86% (133 out of 154) of the responding data manager listed "Lack of sustainable hardware, software or support of computer environment may make the information inaccessible" as high threat to preservation (rating it very important or important). The report on the Warwick Workshop (see [16] ) from 2005 first recommendation for 10 years time is developing massively scalable storage solutions with standard interfaces.

Recommendations on addressing the gaps identified here are discussed in chapter 4.

## 4 RECOMMENDATIONS

Based on the analysis of scalability requirements and the gaps identified, the following recommendations are suggested:

1. **Scalability should be planned as an integral part of the tools and services development.**

There are many aspects of scalability concerning the size and complexity of the data as well as operational requirements from the system and all of them should be considered.

Adding scalability support is not a point action and evolving existing systems to support scalability touches many aspects of the system. Proper planning and resource allocation is required in order to add scalability. Based on the observation that the increase in scalability demands is often underestimated, care should be taken while planning to allow for proper margins of growth.

This recommendation is very much in line with the view of the EU research projects described in section 2.2.3, that either take scalability as part of the overall design (ENSURE, SCAPE) or took at all aspects of a toolset to redesign for scalability (SCIDIP-ES).

Because the recommendation is to look at the full set of scalability requirements with the other requirements of the tool or service, we do not see a place for recommendations for the use of specific tools or techniques in general.

2. **Whenever possible, the use of commercially available systems, tools, and services, instead of in house developed solutions should be considered.**

In many cases, organizations prefer small in house solutions to using external ones. For smaller organizations in many cases it is hard to evolve those solutions to support growing scalability (and other) emerging requirements. When alternatives are available they should be considered, taking into account that preservation systems are intended for long term activity and scalability requirements only grow.

3. **In order to benefit from economies of scale in storage, the use of cloud storage should be considered.** Most organizations equate cloud storage with public clouds. Cloud storage delivery models can be flexible and allow for sharing in controlled ways, as described in section 2.3.3.

The advantages of cloud storage include scalability, elasticity and availability, and enable simple access to data from any location and any device. They usually have a well defined service level agreement (SLA) specifying what the system's properties are. This makes it simpler to build more scalable preservation systems. Another advantage is the pay as you go cost model that makes the investment more proportional to the amount of data.

There are several factors that make organizations reluctant to move data to the cloud, not only public clouds but other could delivery models as well. Among those factors are: security and specifically isolation of tenants in a shared environment; lack of control over the data; concerns over vendor lock-in; as well as others. Those legitimate concerns should be answered by the cloud providers, either via SLAs or by developing better mechanisms in cloud offerings.

In addition, most clouds today do not have support of the full preservation process. Explicit support for preservation should be added to the cloud, either as part of the cloud or as a service above the cloud to make them more usable for preservation.

4. **Consider preservation as part of the overall organization's activity.**

Many of the scalability issues are better handled when the overall organizational activity and forecasts are considered and not as part of an isolated preservation system. This will make the scalability requirements from the preservation part of the overall scalability requirements.

The recommendations above should be considered as part of the APARSEN Virtual Centre of Excellence. The scalability recommendations should be combined with other APARSEN results. Obviously, the other usability areas that include tools, standards and interoperability, are related to scalability. In addition, the work on storage solutions is closely related, as an important part of the

scalability gap is in the storage solutions used. Another related area is cost and cost modelling as any attempt to improve scalability has cost implication. The combined results can help in addressing the scalability gap we are facing and better direct the emergence of preservation systems of the future.

## 5 ANNEXES

### 5.1 QUESTIONNAIRE

APARSEN is a Network of Excellence funded by the EU (2011-2014) with the goal of overcoming the fragmentation of the research and of the development in the digital preservation area by bringing together major European players. Among other activities, APARSEN is investigating the storage solutions currently used in digital archives and preservation repositories, in order to compare the different approaches and technical solutions, and their fitness to preserve digital objects in terms of cost, legal, policy and business needs. Moreover, the project is also addressing analyzing the scalability problems that arise from the increase in number and size of preserved digital objects, that many repositories are currently experiencing.

This questionnaire has been produced as part of this effort, in order to analyze the current practices and to gather comments and suggestions from all the involved parties, based on their individual experiences. The questionnaire is being sent to a number of repositories across Europe, having a variety of profiles that range from large repositories of scientific data, to national libraries and national archives, up to minor digital libraries and repositories run by small organizations with a limited budget. The results of the analysis we plan to perform by investigating such a large audience will hopefully provide an interesting feedback to the digital preservation community, and will allow sharing different experiences.

We do hope that you may give your contribution and cooperate with us by filling in this questionnaire, and we will be grateful to you for doing so. To minimize your effort we have tried to keep the number of questions within reasonable limits, and we have designed the form to speed up the fill in process, by allowing most answers to be given just through checkboxes.

The questionnaire is structured in five sections:

1. Basic information about the institution and the repository
2. Types of digital objects
3. In house storage solutions
4. Outsourced storage solutions
5. Moving to an outsourced storage solution

Sections 1 and 2 are appropriate for any kind of repository, and we expect everyone to complete them. Sections 3 and 4, instead, may be alternative, at least for those repositories relying on a single kind of storage solution (in house or outsourced). As for section 5, due to the relevance of the topic, we will really appreciate if you can answer these questions as well, even if you have currently no intention to move to an outsourced storage solution.

On the other hand, since you may have several types of digital objects and several storage solutions within your repository, additional forms are provided to you for these sections, if you need to fill them. Though almost all questions can be reasonably well understood without the need of specific instructions, we nevertheless provide you with a full set of instructions, that you will find enclosed in the distribution bundle. You do not have, of course, to read them, but just to refer to them, if you need to. The instructions include also a list of selected terms from the Digital Preservation Glossary that is being produced by the APARSEN project.

For any further question, please do not hesitate to ask for assistance the person who has contacted you to send you the questionnaire, or send an e-mail directly to: [salza@dis.uniroma1.it](mailto:salza@dis.uniroma1.it).

## 1. BASIC INFORMATION ABOUT THE INSTITUTION AND THE REPOSITORY

### 1.1. Name of institution:

### 1.2. Contact person

Last name:

Name:

e-mail:

### 1.3. Description (the institution and the repository)

### 1.4. Mission

### 1.5. Ownership:

☐ public

☐ private

### 1.6. Legal regulations (specific national and international regulations to comply with)

### 1.7. Experience: how long have you been preserving digital objects?

- ☐ less than 5 years
- ☐ 5 to 10 years
- ☐ 10 to 20 years
- ☐ more than 20 years

### 1.8. Total number of digital objects preserved:

- ☐ less than 100.000
- ☐ 100.000 to 1.000.000 more
- ☐ than 1.000.000

### 1.9. Total size of the preserved objects:

- ☐ less than 100 TB
- ☐ 100 TB to 1 PB (1000 TB)
- ☐ more than 1 PB

### 1.10. Annual quantity (number of new objects per year):

- ☐ less than 10.000
- ☐ 10.000 to 100.000 more
- ☐ than 100.000

**1.11. Annual size (total size of new objects preserved every year)**

- ☐ less than 1 TB
- ☐ 1 TB to 100 TB
- ☐ more than 100 TB

**1.12. Expected increase in the number of preserved digital objects (refers to 1.8)**

- Expected increase in the next 5 years less than 100%
  - ☐ from 100% to 500%
  - ☐ more than 500%
- Expected increase in the next 10 years less than 100%
  - ☐ from 100% to 500%
  - ☐ more than 1000%

**1.13. Expected increase in the total size of preserved objects (refers to 1.9)**

- Expected increase in the next 5 years less than 100%
  - ☐ from 100% to 500%
  - ☐ more than 500%
- Expected increase in the next 10 years less than 100%
  - ☐ from 100% to 500%
  - ☐ more than 1000%

**1.14. Policy**

- Is there a declared policy on storage of data in place?
  - ☐ Yes
  - ☐ No
- Public link to the policy:
- What are the key criteria of the policy? (*check one or several*)
  - ☐ Regular integrity checks?
  - ☐ Determinations for controlling / monitoring and redundancy
  - ☐ Established and tested workflow for data
  - ☐ recovery? Other (please specify)

**1.15. Additional comments**

## 2. TYPE OF DIGITAL OBJECTS

(if appropriate fill a separate form for each object type)

### 2.1. Description

### 2.2. Producers and Designated Community

### 2.3. Ownership of digital rights:

- ☐ internal  
☐ external

### 2.4. Access policy:

- a) ☐ unlimited ☐ registered users only  
b) ☐ from anywhere ☐ only from specific locations  
c) Are there specific requirements to be met on access time?  
☐ Yes ☐ No

### 2.5. Format(s)

### 2.6. Versions (number of distinct versions preserved for each digital object):

- ☐ 1  
☐ 1 to 5  
☐ more than 5

### 2.7. Size of metadata (ratio between the size of metadata and the size of data)

- ☐ less than 5%  
☐ 5% to 10%  
☐ more than 10%

### 2.8. Connections between objects (are there objects groups that are often accessed together?)

- ☐ Yes  
☐ No

### 2.9. Annual quantity (number of new objects of this type per year):

- ☐ less than 10.000  
☐ 10.000 to 100.000 more  
☐ than 100.000



**2.10. Total number of objects of this type preserved in the repository:**

- ☐ less than 100.000
- ☐ 100.000 to 1.000.000 more
- ☐ than 1.000.000

**2.11. Total size of the of objects of this type preserved in the repository:**

- ☐ less than 100 TB
- ☐ 100 TB to 1 PB
- ☐ more than 1 PB

**2.12. Average access rate (average number of accesses per object per year):**

- ☐ less than 1
- ☐ 1 to 100
- ☐ more than 100

**2.13. Retention period (how long objects of this type must be preserved):**

- ☐ up to 5 years
- ☐ 5 to 10 years
- ☐ more than 10 years

**2.14. Storage levels used for this type of object (referring to sect. 3 and 4)**

- Storage for access
  - ☐ RAM/Solid state
  - ☐ Hard disk
  - ☐ RAID or similar
  - ☐ Tape
  - ☐ WORM disk
  - ☐ CD or DVD
- Storage for preservation (if different)
  - ☐ RAM/Solid state
  - ☐ Hard disk
  - ☐ RAID or similar
  - ☐ Tape
  - ☐ WORM disk
  - ☐ CD or DVD
- Storage for further backup copies
  - ☐ (if any) RAM/Solid state
  - ☐ Hard disk
  - ☐ RAID or similar
  - ☐ Tape

- ☐ WORM disk  
☐ CD or DVD

### 2.15. Additional comments

## 3. IN HOUSE STORAGE SOLUTIONS

(if necessary fill a separate form for each different storage solution you are currently managing)

### 3.1. Description

### 3.2. Purpose:

- ☐ access only  
☐ preservation only  
☐ both

### 3.3. Technology:

- ☐ RAM/Solid state  
☐ Hard disk  
☐ RAID or similar  
☐ Tape  
☐ WORM disk  
☐ CD or DVD

### 3.4. Experience: how long have you been using this storage solution?

- ☐ less than 5 year  
☐ 5 to 10 years  
☐ more than 10 years

### 3.5. Redundancy (not including backups)

- **Level:** ☐ none ☐ RAID1 ☐ RAID5 ☐ other  
– **Mode:** ☐ device ☐ local ☐ geographical (more than 50 km)

### 3.6. At which intervals are backups performed for this level of storage?

- ☐ 1 to 7 days  
☐ 7 to 30 days  
☐ 30 to 180 days  
☐ more than 180 days  
☐ never

**3.7. Staff (number of employees devoted to the management of this storage system)**

- ☐ less than 5  
☐ 5 to 10  
☐ more than 10

**3.8. Cost estimate (total cost per TB/year): €****3.9. Which controls are you regularly performing, and at which intervals, to check the integrity of stored data?****3.10. Perceived weakness: why do you think the solution is not completely satisfactory? How could it be improved?****3.11. Time horizon: how long do you plan to reasonably carry on with this solution before replacing the devices and/or moving to a different storage solution?**

- ☐ less than 1 year  
☐ 1 to 5 years  
☐ more than 5 years

**3.12. Risk assessment ( with regard to this storage level)**

- *Reliability*, i.e. probability of losing (part) of your stored data:

- ☐ I cannot provide any reliable estimate  
☐ Yes, my estimate of system reliability is (express as power of 10) :

- *Availability*, fraction of time during which access to the preserved objects is

- ☐ granted: I cannot provide any reliable estimate  
☐ Yes, my estimate of system availability is (express as power of 10):

**3.13. Additional comments**

#### 4. OUTSOURCED STORAGE SOLUTIONS

*(if necessary fill in a separate form for each different storage service you are currently purchasing)*

##### 4.1. Description

##### 4.2. Name of the provider:

##### 4.3. Classification:

- ☐ public cloud
- ☐ private shared cloud: jointly managed by a consortium of several
- ☐ repositories redundant cloud: data stored in several clouds, managed
- by different vendors
- ☐ other (specify)

##### 4.4. Experience: how long have you been using this storage solution?

- ☐ less than 1 year
- ☐ 1 to 5 years
- ☐ more than 5 years

##### 4.5. Redundancy (as declared by the provider)

##### 4.6. Terms of agreement (what does your provider guarantee you, and in which terms):

###### Access

###### Security

###### Persistence

##### 4.7. Storage cost (TB/year): €

##### 4.8. Access cost. a) per access: € b) per MB transferred: €

##### 4.9. Risk assessment ( with regard to this storage level)

- *Reliability*, i.e. probability of losing (part) of your stored data:

- ☐ I cannot provide any reliable estimate
- ☐ Yes, my estimate of system reliability is (express as power of 10)

- *Availability*, fraction of time during which access to the preserved objects is

- ☐ granted: I cannot provide any reliable estimate  
☐ Yes, my estimate of system availability is (express as power of 10)

## 5. MOVING TO AN OUTSOURCED STORAGE SOLUTION

### 5.1. Which outsourcing scheme are you willing to consider, if any? (Check one or several)

- ☐ public cloud  
☐ private shared cloud: jointly managed by a consortium of several repositories  
☐ redundant cloud: data stored in several clouds, managed by different vendors  
☐ other (specify)

### 5.2. Motivations: why could you consider moving to a different storage solution? (Check one or several)

- ☐ improve cost  
☐ improve reliability  
☐ improve availability  
☐ improve scalability  
☐ other (specify)

### 5.3. Fears: why would you refrain from moving to an outsourced storage solution? (Check one or several)

- ☐ lack of control  
☐ fear of unauthorized access to data  
☐ fear of loss of data 1: technical problems fear of  
☐ loss of data 2: provider bankruptcy  
☐ fear of loss of data 3: becoming unable to pay the fees because of financial or cash  
☐ problems vendor lock-in: problems in moving data to a different provider

### 5.4. Additional comments

## 5.2 ANNEX II: QUESTIONNAIRE INSTRUCTIONS

### 1. BASIC INFORMATION ABOUT THE INSTITUTION AND THE REPOSITORY

#### 1.1 Name of the institution

*Name of the institution that manages the repository.*

#### 1.2 Contact person

*Person that should be contacted, if necessary, for further information.*

#### 1.3 Description

*Short description of the repository and of the institution that manages it.*

#### 1.4 Mission

*Mission that the institution has been given as a whole, and specifically with reference to the repository.*

#### 1.5 Ownership

*Is your institution a private or public organization?*

#### 1.6 Legal regulations

*Are there any specific legal regulations (national or international) that apply to the objects preserved in the repository? For instance, special regulations apply, in general, to medical records.*

#### 1.7 Experience

*How long have you been preserving digital objects, and providing this a reliable service?*

#### 1.8 Total number of preserved digital objects

*By digital object we mean every single object that is individually preserved in the repository, and that can be individually retrieved.*

#### 1.9 Total size of preserved digital objects

*Total size does not include backup copies or replications, created or managed within the repository for whatever purpose. It just accounts for the original size of the preserved objects.*

#### 1.10 Annual quantity

*Estimate of the average number of new objects that are ingested by the repository every year.*

#### 1.11 Annual size

*Estimate of the average yearly increase in the total size of preserved digital objects (see 1.9).*

#### 1.12 Expected increase in the number of preserved digital objects

*How much do you expect the total number of preserved object (see 1.8) to increase in the next 5 and 10 years*

#### 1.13 Expected increase in the total size of preserved objects

*How much do you expect the total size of preserved object (see 1.9) to increase in the next 5 and 10 years*

#### 1.14 Policy

*Is there a formal and declared policy on storage of data the repository must conform with. Such a policy should be stated in a formal document. Specify the key criteria according to the list or add additional ones.*

#### 1.15 Additional comments

*Add any further information you may like to provide.*

## 2. TYPES OF DIGITAL OBJECTS

*If your repository preserves several kinds of digital objects that are substantially different, you may want to fill in several copies of this section. You may find in the distribution bundle a separate file with a copy of this section. You may duplicate it as you need.*

### 2.1 Description

*Short description of this type of digital objects, to clearly identify what they are.*

### 2.2 Producers and Designated Community

*Specify the organizations that provide the information to be preserved, and community of users that is interested in preserving the digital objects, and in being able to continue accessing them (see glossary).*

### 2.3 Ownership of digital rights

*Who owns the digital rights on the preserved objects? Is it the institution managing the repository or not?*

### 2.4 Access scheme

*a) Is access limited to registered users or registered users?*

*b) Is access limited to local users?*

*c) Are there specific requirements on access time (e.g. maximum access time)?*

### 2.5 Formats

*Format of the digital objects and/or of their components if they have a complex structure.*

### 2.6 Versions

*Are several versions of the same object typically preserved in the repository as a single object?*

### 2.7 Size of metadata

*Quantify the impact of metadata on the total amount of stored data.*

### 2.8 Connections between objects

*Are there connections between objects, such that groups of objects are preserved and often accessed together?*

### 2.9 Annual quantity

*Average number of new objects of this type ingested by the repository every year.*

### 2.10 Total number of objects preserved

*Total number of objects of this type currently preserved in the repository.*

### 2.11 Total size of the preserved objects

*Total size of the objects of this type currently preserved in the repository.*

### 2.12 Average access rate

*Average number of accesses per year and per preserved object.*

### 2.13 Retention period

*How long must the objects of this type be preserved, according to the regulations the repository must comply with.*

### 2.14 Storage levels

*Referring to the storage solutions described in section 3 and 4 of the questionnaire, specify which ones are used for access copies, for preservation copies and for backups.*

### 2.15 Additional comments

*Add any further information you may like to provide.*

### 3. IN HOUSE STORAGE SOLUTIONS

*In house solutions are opposed to outsourced storage solutions, as for instance service provided by storage farms or by the cloud. If your repository manages several kinds of storage solutions that are substantially different, you may want to fill in several copies of this section. You may find in the distribution bundle a separate file with a copy of this section. You may duplicate it as you need.*

#### 3.1 Description

*Short description of this type storage system the form refers to.*

#### 3.2 Purpose

*Are the copies of digital objects stored in this storage system used only for access, only for preservation or for both?*

#### 3.3 Technology

*Which technology exploits this storage system? If appropriate, if for instance several layers of storage are included, you may check several.*

#### 3.4 Experience

*How long is your experience in using this specific solution?*

#### 3.5 Redundancy

*Some storage solutions involve storing several replications of the same object. This may be achieved by replicating the data within a single device, e. g. a RAID 5 NAS, or by storing different copies in several devices in the same location, or in different locations.*

#### 3.6 Backups

*Specify at which intervals are backup performed. By backup we mean generating additional copies of the preserved objects to be stored and preserved elsewhere.*

#### 3.7 Staff

*How many of your employees are entirely devoted to manage this storage system?*

#### 3.8 Cost

*Provide even a rough estimate of the Total Ownership Cost of this storage system. Includes everything: purchase, maintenance, personnel, rental cost etc.*

#### 3.9 Controls

*Specify which kind of controls (if any) are regularly scheduled to check the integrity of the data stored within this storage system.*

#### 3.10 Perceived weakness

*Do you rate this solution as completely satisfactory, or you perceive some weakness in it?*

#### 3.11 Time horizon

*How long are you planning to continue using this solution, before obsolescence or other problems may force you to move to a different one.*

#### 3.12 Risk assessment

*The probability of losing the data refers to this level of storage, i.e. means you become unable to access them within this storage system. This does not means that you have lost your data, since you may still be able to access redundant copies or backups stored in different systems. The same goes with availability.*

#### 3.13 Additional comments

*Add any further information you may like to provide.*



## 4. OUTSORCED STORAGE SOLUTIONS

*If your repository relies on several storage services, you may want to fill in several copies of this section. You may find in the distribution bundle a separate file with a copy of this section. You may duplicate it as you need.*

### 4.1 Description

*Specify the kind of service and how it is provided*

### 4.2 Name of the provider

*Name of the vendor or of the organization that is providing you the service?*

### 4.3 Classification

*Check the appropriate type or provide your own, by giving a concise definition of it.*

### 4.4 Experience

*How long have been using this storage service?*

### 4.5 Redundancy

*Does the provider specify which kind of redundancy is implemented to guarantee the persistence of the stored data? Sometimes it does.*

### 4.6 Terms of agreement

*Persistence, access and security (confidentiality of the data) are the three main issues. What your provider guarantees you and in which terms.*

### 4.7 Storage cost

*There may be different billing schemes. Some include everything, some bill access independently from storage. This should be the basic cost, including storage, but not necessarily access.*

### 4.8 Access cost

*According to different schemes, one may pay per access and/or per MB.*

### 4.9 Risk assessment

*The probability of losing the data refers only to this level of storage, i.e. means you become unable to access them through this storage service. This does not mean that you have lost your data, since you may still be able to access redundant copies or backups stored in different systems. The same goes with availability.*

## 5. MOVING TO AN OUTSORCED STORAGE SOLUTIONS

*By outsourced storage solution we mean buying storage as a service commercially sold by a third party, instead of setting up and managing directly a private storage. Even if you have currently no intention to move to such storage solution, it may still be very interesting that you fill in this section of the questionnaire, at least to specify which outsourced storage solutions you have considered (or you are considering) and which could be, in your opinion, the pro and cons of such a move.*

### 5.1 Which outsourcing schemes are you willing to consider?

*By public cloud we mean networked [online storage](#) commercially sold as a service by a vendor. Instead, by private cloud we mean a similar service set up by a consortium, with access limited to the member of the consortium. Setting up a private cloud allows typically to meet specifications that are not met by public clouds (at least in the opinion of consortium members). Finally by redundant cloud we mean a scheme where data are stored in several public clouds, managed by different vendors, so providing further redundancy, and potentially avoiding vendor lock-in.*

### 5.2 Motivations

*Specify which advantages you may expect from moving to an outsourced storage solution. The most likely motivations are listed, but, if you can see a different one, please indicate it.*

### **5.3 Fears**

*Specify which disadvantages you may expect from moving to an outsourced storage solution. The most likely motivations are listed, but, if you can see a different one, please indicate it.*

### **5.4 Additional comments**

*Add any further comment you may like to provide.*

## 5.3 ANNEX III: QUESTIONNAIRE GLOSSARY

### Access Policy

- Written statement authorized by the repository management that describes the approach to be taken by the repository for providing access to objects accessioned into the repository. The Access Policy may distinguish between different types of access rights, for example between system administrators, Designated Communities and general users. OAIS
- A set of rules and guidelines that determine how the institution's collections, services, products and databases are accessed. ANZ

### Archive

- An organization that intends to preserve information for access and use by a Designated Community. OAIS
- An organisation (or part of an organisation) responsible for appraising, acquiring, preserving and making material available. ANZ

**Designated Community:** An identified group of potential Consumers who should be able to understand a particular set of information. The Designated Community may be composed of multiple user communities. A Designated Community is defined by the archive and this definition may change over time. OAIS

**Digital Archive:** A repository for the long-term maintenance of digital resources and making them available. ANZ

### Digital Object

- An object composed of a set of bit sequences. OAIS
- Digital objects are units of content managed by a digital archive. Digital objects have as one of their attributes an Identifier. They can be seen to be the atomic level of content. Smaller units can be contained within them, but the repository manages the digital object in a singular fashion. ANZ

### Digital Preservation

- Refers to the series of managed activities necessary to ensure continued access to digital materials for as long as necessary. Digital preservation is defined very broadly for the purposes of this study and refers to all of the actions required to maintain access to digital materials beyond the limits of media failure or technological change. Those materials may be records created during the day-to-day business of an organisation; "born-digital" materials created for a specific purpose (e.g. teaching resources); or the products of digitisation projects. This handbook specifically excludes the potential use of digital technology to preserve the original artefacts through digitisation. DPC
- Principles, practices, methods, strategies and managed activities that ensure long term preservation for continued access to digital materials for as long as necessary. ANZ

**Digital Rights Management:** An umbrella term referring to any of several technical methods used to control or restrict the use of digital content. ANZ

### Format

- Specific, pre-established structure for the organisation of a file or bitstream. ANZ
- "The physical medium in which information is recorded or carried, e.g. paper files, computer printout, photographs, microfilm, [electronic] records, plans, cards, volumes, etc. A selection of descriptive elements set out in a prescribed manner and sequence so that the resulting description will be standardised for all types of records." ANZ

### Metadata

- Data about other data. OAIS
- Information which describes significant aspects of a resource. Most discussion to date has tended to emphasise metadata for the purposes of resource discovery. The emphasis in this handbook is on what metadata are required successfully to manage and preserve digital materials over time and which will assist in ensuring essential contextual, historical, and technical information are preserved along with the digital object. DPC
- Structured information that describes and/or allows users to find, manage, control, understand or preserve information over time. Source: ICA req OR "Data describing context, content and structure of records and their management through time." (ISO 15489-2001) ANZ

**Preservation:** The processes and operations in ensuring the technical and intellectual survival of objects through time. Source: derived from Continuum ANZ

**Producer:** The role played by those persons, or client systems, who provide the information to be preserved. This can include other OAISs or internal OAIS persons or systems. OAIS, ANZ

**Repository:** "The building or room, or part thereof, set aside for the storage of archives and/or intermediate records. Archival repositories are often constructed to meet specific environmental standards designed to ensure the longevity of the records." ANZ

**Retention Period:** "The period of time, usually based on an estimate of the frequency of current and future use, and taking into account statutory and regulatory provisions, that records need to be retained before their final disposal. Sometimes used to indicate the length of time records are to be retained in offices before being transferred to intermediate storage." (KA, p.479). The retention period usually commences from the time of the disposal trigger. ANZ

**Version:** An attribute of an AIP whose information content has undergone a transformation on a source AIP and is a candidate to replace the source AIP. OAIS

## Glossary References

*OAIS - Open Archival Information System - Reference Model - ISO 14721:2012*

<http://www.alliancepermanentaccess.org/index.php/knowledge-base/member-resources/digital-preservation-glossary/>

*DPC - Digital Preservation Coalition*

<http://www.dpconline.org/advice/preservationhandbook/introduction/definitions-and-concepts>

*ANZ - Archives of New Zealand*

## 6 REFERENCES

- [1] ISO 14721:2003, "Space Data and Information Transfer Systems - Open Archival Information System - Reference Model", Edition 1, February 2003.
- [2] CCSDS 650.0-B-1, "Reference Model for an Open Archival Information System (OAIS)" – CCSDS Blue Book, Issue 1, January 2002.
- [3] "THE DIGITAL UNIVERSE IN 2020", December 2012, [www.emc.com/collateral/analyst.../idc-the-digital-universe-in-2020.pdf](http://www.emc.com/collateral/analyst.../idc-the-digital-universe-in-2020.pdf)
- [4] "2010 Digital Universe Study: A Digital Universe Decade – Are You Ready?", 4-26-2010, [http://gigaom.files.wordpress.com/2010/05/2010-digital-universe-iview\\_5-4-10.pdf](http://gigaom.files.wordpress.com/2010/05/2010-digital-universe-iview_5-4-10.pdf)
- [5] "The Expanding Digital Universe", March 2007, <http://www.emc.com/collateral/analyst-reports/expanding-digital-idc-white-paper.pdf>
- [6] Amazon Web Services Blog, April 05, 2012, <http://aws.typepad.com/aws/2012/04/amazon-s3-905-billion-objects-and-650000-requestssecond.html>
- [7] "SCIDIP-ES: D32.1 Generic Services/Toolkits and Scalability Research Report and Plan", <http://www.scidip-es.eu/scidip-es/deliverables/SCIDIP-ES-DEL-WP32-D32.1.pdf>
- [8] Mell, P. and Grance, T., "The NIST definition of cloud computing", National Institute of Standards and Technology, 53(6), 2009.
- [9] McLean, B., "DuraCloud: Federated Repositories and Cyberinfrastructure Open technologies and services for managing durable data in the cloud", 2009.
- [10] Iglesias, E. and Meesangnil, W., "Using Amazon S3 in Digital Preservation in a mid sized academic library: A case study of CCSU ERIS digital archive system", *Code{4}lib Journal*, 12/2010.
- [11] Cloud Data Management Interface, Version 1.0.1, SNIA Technical Position. September 2011. <http://cdmi.sniacloud.com>.
- [12] David Rosenthal. Preservation in the cloud. In Preservation in the Cloud. Library of Congress, September 2009.
- [13] Long-term Preservation Storage: OCLC Digital Archive versus Amazon S3. <http://dltj.org/article/oclc-digital-archive-vs-amazon-s3/>.
- [14] Blevins, R., Patel, I., O'Sullivan, J., Hunter, I., and Sharpe, R., "Studies on the scalability of web preservation", to be presented in the 10th International Conference on Preservation of Digital Objects (iPRES 2013) 2-6 September 2013, Lisbon, Portugal
- [15] "Insight into digital preservation of research output in Europe", PARSE.Insight Deliverable: D3.4, 9 December 2009, [http://www.parse-insight.eu/downloads/PARSE-Insight\\_D3-4\\_SurveyReport\\_final\\_hq.pdf](http://www.parse-insight.eu/downloads/PARSE-Insight_D3-4_SurveyReport_final_hq.pdf)
- [16] "Digital Curation and Preservation: Defining the research agenda for the next decade", Report of the Warwick Workshop - 7 & 8 November 2005, Report Co-ordinators: David Giarretta, Heather Weaver
- [17] "The Ability to Preserve a Large Volume of Digital Assets - A Scaling Proof of Concept", Ex Libris, 2010, <http://www.exlibrisgroup.com/files/Products/Preservation/RosettaScalingProofofConcept.pdf>