



An Enterprise Strategy for Low Cost Digital Storage

A small library cannot become a big one just by building more cheap shelving. It requires shelving appropriate for different resources. ... and it requires a catalogue system to access content by its attributes, and acquisition and collection weeding practices to manage the value of the collection.

The demand for storage continues to grow. While part of this growth demands high performance storage for effective online applications, a growing portion of this storage has different characteristics that suggest lower cost storage alternatives. To meet storage demands, an enterprise storage infrastructure must include technology with different cost/performance characteristics.

Managing cost also includes comprehensive data management practices and tools to ensure information is stored efficiently and securely. A comprehensive Digital Storage Plan is the key.

- **Data management strategy** with institutional information management policies and guidelines aware of the content offered and stored.
- **Storage infrastructure strategy** with alternative storage services at pricing sensitive to operating costs, unit/user ability to pay, and institutional resource management strategies.

As a starting point, this paper considers classes of storage requirements to identify needed storage components that are more cost effective than the current resource. The volume of much of the information currently stored could be reduced considerably and stored more efficiently, reducing infrastructure costs while also reducing DR backup frequency and capacity considerably.

Since there are a wide array of lower cost storage technologies with different characteristics, it is recommended that the vendor community be engaged in developing a Digital Storage Plan.

Given a large portion of the information stored can be supported with lower performance infrastructure, CCS should prioritize the implementation of a low-cost, open source, commodity-driven storage system in parallel with the Network Appliance. A rapidly evolving technology niche, it requires a strategy of DIY and continuing evolution. Such a strategy can be supported by a strategic vendor partner.

An Enterprise Strategy for Low Cost Digital Storage

Executive Summary

Request

This investigation reviews the spectrum of needs for low cost storage services identified in the community and considers how CCS could deliver low cost enterprise storage.

Findings

While the “standard” request of CCS is for larger amounts of lower-cost storage, the analysis of requests from various users and groups in the University identify two interrelated issues

Growing Demand for Storage

The volume of digital content is growing rapidly because of interrelated issues.

- Rapid production all kinds digital materials.
- An understanding that old content is too valuable to throw away.
 - compounded by the lack of institutional policies that define responsibilities.
- Volume of stored data is too large for individuals to manage effectively.

Storage is consumed with unnecessary or obsolete content. **CCS should provide leadership in the development of Digital Content ownership/stewardship and retention policies.**

Classification of Storage Requirements

Based on reported studies and a small sample of University persons looking for additional storage, an estimate of the profile of storage requirements is described in **Table 1**. This classification implies that the enterprise storage infrastructure strategy should be a mix of cost/performance technologies, biased primarily towards low cost over high performance storage.

CCS requires better storage content analysis tools and processes to determine the University’s actual mix of storage requirements.

Demand for Low Cost Alternatives

The demand for low cost alternatives stems from several issues that result in users looking for costs in line with commodity workstation hard drives.

- Limited attention to total information resource management costs.
- A funding model that emphasizes initial cost avoidance vs. TCO for decision maker.
- Information management environment promotes high risk alternatives.
- Commodity web-based services promote “storage is free” perception.

Information Management Strategy

While the classification model recommends expansion of low-cost storage capacity, providing only more capacity will only exacerbate volume growth. It will be more cost effective in the longer term to front-end low-cost storage with applications and services that promote effective information management. The starting point is an archive system.

Low Cost Storage

Implementing a low-cost storage capacity is a priority. However,

- A solution independent of the current NetApp infrastructure will add considerable operating cost. **CCS should expand storage based on a clearly defined integrated solution**, rather than adding complexity to the current storage infrastructure. Backup must be considered integral to a storage plan, not an additional function.

- There are a number of integrated storage strategies supported by major storage vendors and integrators. Investigation of alternative components and compatibility is too complex for CCS. These vendors have a better perspective on storage technology, alternatives and the necessary management tools needed for cost effective operation. Before adding additional capacity, **CCS should work with appropriate consultants and vendors to select a manageable storage solution.**
- Storage technology is changing rapidly. CCS must determine how aggressive its storage renewal plan is. Based on examples of some other larger universities, *CCS should focus on an integrated solution from a proven vendor while also exploring leading edge alternatives.* Two such opportunities are
 - **Massive-Array-of –Idle-Disks (MAID)** storage provides high-density, low-energy-use storage for infrequently accessed content such as research data.
 - **Cloud Archive** services provide cost-effective historical record storage, relying on the provider for disaster recovery and technology renewal.

Summary of Recommendations

Strategic Planning

Recommendation 1 that CCS develop a model plan for an Integrated Storage Infrastructure strategy based on a current industry model.

Recommendation 2 that CCS encourage and support the development of institutional information management policies that clarify responsibilities for digital content management.

Tactical Direction

Recommendation 3 ... that CCS do a storage content analysis of enterprise and departmental storage to determine actual mix of storage classes.

Recommendation 4 ... that CCS immediately begin an RFI process that solicits recommendations on a suitable storage infrastructure plan and qualifies strategic vendors prior to an RFP for specific components.

Recommendation 5 ... that, as an interim step, CCS purchase a lower cost commodity SAN to extend CFS storage, implemented clearly as lower-performance, additional-drive storage with long but less frequent DR cycles. This solution would be replaced by the outcome of Rec. 4.

Recommendation 6 ... that CCS develop an information archive service, using existing storage capacity to begin to engage the community in information management and to off-load and re-direct storage requirements with lower access performance needs.

- A Research data store might be a priority special archive.

Recommendation 7 ... that CCS initiate a project to investigate the opportunity/complexity of introducing lower cost open source and commodity components.

Recommendation 8 ... that CCS implement a small scale “Permanent Historical Record” archive with a cloud provider to evaluate and demonstrate this opportunity.

Table 1: An estimate of a storage requirements profile

<u>Class</u>	<u>Examples</u>	<u>Estimated Portion</u> *	<u>Access Performance</u>	<u>Rate of change</u>	<u>Possible DR plan</u>
Transactional <i>accessed through online processing systems; computation for research</i>	ERP systems HPC computing Web service backend	5%	High	Continual	Frequent database-level backups/snapshots
Online Dynamic Content <i>Accessed through workstation file systems (CFS, NFS)</i>	Document storage CFS General NFS	10%	Medium	Frequent; timely	Weekly Full Daily Incremental Keep: 1 semester
Online Static Content <i>Accessed through web-based applications</i>	Completed Documents. Institutional reports. Web pages / Web CMS.	20%	Medium – Low (cached for web performance)	Limited; not timely	Monthly Full Weekly Incremental Keep: 1 year
Infrequent Use but Accessible <i>Information used on an as-needed basis, not requiring performance access</i>	Older document versions accessed monthly/yearly. Active Research Data analysed on cloud or shared HPC services.	25%	Low (similar to Internet access)	Read-only after create;	Create ILM metadata; Weekly Catalogue audit; Weed monthly based on ILM metadata; DR Backup on create; delete on weed;
Archive <i>Information stored mainly for long term preservation; accessed vary infrequently</i>	Research Data Preserved documents Old Financial Records Old Student Records	40%	Low (request/retrieve in minutes)	Read-only; deletion protected.	Create ILM metadata; Weekly Catalogue audit; Weed annually based on ILM metadata; Self-contained DR plan (e.g. mirrored sites)

* The Estimated Portion of total storage demand is obtained from 200-2008 studies and hearsay Guelph information. The general industry understanding is that total ERP Transactional content is not growing rapidly because of year-end type paring. However, since other content is being kept longer the significant growth is in the infrequent and archive areas. A leading consultant recently commented that most corporations now have more than 50% of their storage consumed by archive-class information.

An Enterprise Strategy for Low Cost Digital Storage

Request

This investigation reviews the spectrum of needs for low cost storage services identified in the community and considers how CCS could deliver low cost enterprise storage.

Findings

Storage Demand

There are two issues associated with storage demand.

1. Growing Demand for Storage

Faced with growing volume of content, many users deal lack of storage by buying larger hard drives and avoiding the task of managing older/obsolete content. With enterprise provided storage, the same approach is used – request more storage space.

The volume of digital content is growing rapidly because of a layered

- **Rapid production** - Almost all information is being created in digital form. Multimedia content is very easy to create. Administrative, learning, and research all create information rapidly.
- **Old content is valuable** - Many users avoid throwing away old content. Research data usually has high history value, but value may deteriorate if related information becomes dissociated.
 - This problem is compounded by the lack of institutional/departmental information management policies that results in users being afraid to delete old content.
- **Duplication** - Besides the obvious duplication of systems files and application software, users are duplicating stored information.
 - Without content management versioning, older versions are left in storage.
 - Much information needs to be shared; files are copied to the storage of others.
 - As a mitigation of high risk local storage, user create their own backups.
- **Type of Content** - Multimedia content and large research datasets consume large amounts of space and have different storage and preservation requirements. There is a growing need to store “packages” rather than individual files. (e.g. a collection of research pictures, or metadata descriptions with data files, ...)

2. Demand for Low Cost Alternatives

Individual users (and research groups) are recognizing some of the inherent operational effort and risks associated with storage of large volumes of information. (Note “large” is a relative term influenced by a person’s comfort with digital technology and with the type of content. Managing family photos is different than managing a large number of aging research spreadsheets.)

Matching content value to storage cost

Many administrative users now understand value in managed enterprise storage. However, without clear institutional policies on content stewardship and retention expectations, users attach a low value on much of their content. They look for a low-cost enterprise-managed service that matches this low content value.

Similarly as research projects are completed, researchers place a limited value on their original data. Historically, they have just placed a copy of the data on their office shelves! However, they now recognize the value of maintaining the data in digital form. With no on-going funding they are looking for enterprise provided retention storage.

Funding Storage Costs drives decisions

Individual staff or workgroups seem to be responsible for a large amount of the University's enterprise digital content, including where it is stored and how the content is maintained. There are few institutional content management policies that direct staff on how to manage this content. In many cases, individuals are responsible for making storage procurement decisions (mostly involving desktop/laptop storage). Most of these decisions are based on capital cost of the storage device, not on the total cost of ownership (TCO).

Researchers are faced with the dilemma of retaining their data while research funders and the local research enterprise do not provide the funding for long term storage. Libraries and discipline groups talk about regional/national/discipline repositories. If they exist, most researchers are not aware of funded repositories. The result is they have to rely on small discretionary funding to pay for low cost storage devices.

The demand for low cost alternatives stems from several issues that result in users looking for costs in line with commodity workstation hard drives.

- **Value/Cost** - While users disposing of old content, they attach an apparent low value to historical data so they are looking for the lowest initial cost solution.
- **Funding model** - Individual users, rather than departments, are making storage procurement decisions, electing to minimize initial cost rather than TCO.
- **High Risk** - Users with workstation storage operate at high risk, minimizing DR and data management costs. They do not look at costs associated with future availability.
- **Web-based alternatives** - Many users are aware of free/low-cost web-based storage for targeted applications (Gmail, Picasso, ...). The assumption is that large volume storage is inexpensive.
- **Managing large content volume** - Most users have no or limited, application-specific tools to manage large volume of data.

Classification of Storage Requirements

Various Information Lifecycle Management (ILM) papers from leading vendors and consultants note the importance of ILM in managing storage costs, both by utilizing the more cost effective storage technology for the type of information and by reducing storage of unnecessary and redundant information.

The results of classification of storage varies widely by type of business. Since the University does not audit stored information by any classification, an actual profile is not possible. However, based on reported studies and a small sample of University of personal descriptions from a few staff looking for additional storage, an estimate of the profile of storage requirements is described in **Table 1**.

This classification suggests that the storage infrastructure strategy should be biased toward technologies that favour low cost over high performance. Business Continuity solutions must be integrated with the storage plan, not a separate service.

Application of a classification of information not only applies to the active data storage systems, but has a major impact on Disaster Recovery (DR) Backup capacity and management processes.

Low Cost Storage Alternatives

Storage costs include both infrastructure costs and storage management overhead. In the current model, CCS charges only for infrastructure costs. To accommodate demands for lower cost storage, a variety of storage technologies with difference performance characteristics tied to specific services would be required. Unless a carefully chosen, integrated, storage management platform is also implemented storage administration effort and costs will increase sharply.

As content volume and complexity increase, our enterprise existing storage solutions will not scale well unless CCS can develop a good picture of the types of information stored

Storage costs include both infrastructure costs and storage management overhead. In the current model, CCS charges only for infrastructure costs. To accommodate demands for lower cost storage, a variety of storage technologies with different performance characteristics tied to specific services will be required.

Having different storage platforms increases complexity. Users can't be expected to manage multiple stores. As information content evolves, it changes classification and should be moved to an appropriate storage. Unless a carefully chosen, integrated, storage management platform is also implemented storage administration effort and costs will increase sharply if multiple storage technologies are used.

Storage Platform Alternatives

Storage systems are designed and marketed to optimize cost for a fairly narrow set of performance criteria. To optimize cost across storage needs that have a wide variety of performance criteria, an enterprise strategy falls somewhere between these two extremes.

One Storage System / Multiple Services priced differently

A single storage system configured to meet most performance criteria is more easily managed and is more adaptable to changing requirements. "Virtual" multi-tier services are easy to deliver. Service pricing can be established depending on "market" requirements with a goal for overall cost recovery. Lower prices of low end services must be balanced with higher prices for high end services.

The performance of a single storage platform must meet the highest requirements. *With a storage mix requiring mostly low/mid performance, a single storage solution is not likely economical.*

Multiple Storage Systems / Individual Service Pricing

Selecting individual storage systems matched to individual service needs lowers the capital cost of storage. (Note that capital costs are still repeating because of the rapid evolution of technology.) Prices can be related more closely to costs.

However, unless the different storage systems are integrated through a single vendor strategy or by an overall non-vendor-specific storage management system, the overall operating cost will be higher as much of the content will eventually move across storage services.

It is a mistake to expect the user community to manage content across a variety of systems, even if driven by pricing!

Two platform approach

The approach recommended by some consultants and the somewhat accidental approach of a number of larger higher education institutions is to build two storage systems, one from each end of the requirements spectrum.

1. A high integrated performance storage solution to support enterprise class applications (ERP) and other similar high performance applications. A variety of performance services are delivered.
2. A low-end storage system that meets the needs of applications with large volume, low performance requirements. Several services can be delivered on this storage.

Cloud Storage Services

At the present, "cloud" storage provides similar low-end functionality with slightly more constrained limits (volume/bandwidth). Cloud storage may be appropriate for some services. However, contractual and security concerns limit its use.

Cloud Storage costs are comparable with enterprise solutions. However, most providers have bandwidth charges for accessing content to pay their networking costs. This results in a higher- than-desirable cost for storage of active information data. However, since infrequent access archival storage may be cost effective in the Cloud.

Conclusions

CCS must “know” the Information it stores

Enterprise Storage is one of the key mission critical infrastructures of the University. Everything is being stored in digital form. Active current content, as well as retained historical content is critical to all systems and services of learning, research and administration. Because content stored digitally is fragile, active file backups and Disaster Recovery backups are an integral part of an enterprise storage service.

Enterprise information content has a broad set of characteristics and hence a range of availability, performance, and lifecycle characteristics. Enterprise information is generated and used through a wide variety of applications, each with their own expectations of the storage service.

To store enterprise information reliably, securely, and efficiently, CCS must be able to audit and analyse enterprise content to classify it into manageable similar-characteristics classes in order to identify efficient storage infrastructure and to apply some policies broadly to classes of information.

CCS Leadership towards Enterprise Information Management Policy

The critical issue about large volume storage is how to manage it effectively. While all organization struggle with the problem, commercial organizations primarily develop enterprise policies and processes to handle most of the content. A university context is different. Policies favour personal ownership and privacy. Individual users, or small workgroups, are primarily responsible for managing the digital information they produce and/or accumulate.

The alternative to central content management policy is the delivery of multiple services to the end user that permits them to exercise better content management.

There is no evidence to suggest that end-user managed storage can be cost efficient unless users have effective tools to manage the life cycle of their content and clear distributed information stewardship policies. Such services must assist users to identify content classes and to easily apply content management (information life cycle) policies to either the classes or individual digital objects.

CCS needs a Comprehensive Enterprise Storage Plan

With a better understanding of information stewardship policies, and a better understanding of the performance requirements of information stored, CCS can develop a longer term storage infrastructure strategy.

Storage technology continues to evolve quickly with promising new technologies “in the lab”. It will be difficult for CCS to keep up-to-date on the breadth, and cost, of solutions that will become mainstream in the next few years.

An optimum storage solution is a combination of advanced technologies matched with enterprise priorities and information storage requirements. That matching requires a knowledge base and expertise beyond CCS capacity.

It is important for CCS to develop a close partnership with a strategic storage vendor who can provide leading industry knowledge and design expertise.

Gartner recommends that an RFI process be developed to explore the regional marketplace for capable partners and to build a CCS understanding of the range of solutions alternatives.

The RFI would simply state the UOG storage environment and request a recommended longer term strategy that the vendor could support. Responses should demonstrate the ability to define and deliver on a long term strategy that is aware of future technology advances.

Coupled with an RFP process for solutions to immediate needs the RFI/RFP provides a way to understand vendors and build a strategic relationship.

CCS should broaden its knowledge with immediate needs projects

Lower cost storage solutions requires consideration of lower capital cost components and some DIY implementations, rather than higher-up-front-cost full-featured monolithic storage solutions.

This approach is growing in popularity but it requires a higher operational effort and potentially lower reliability. Implementing such a solution is appropriate for an appropriate class of storage. In-frequently accessed, fairly static, content that does not require high performance during use (e.g. not databases!) would work well in a low cost storage system.

Open Source Platforms

According to Gartner, a key to reducing cost is to avoid high capital costs on components that will be replaced sooner-than-later. One option that is demonstrating success in the higher education sector especially is the choice of open source software as a platform. (*CCS's choice of Zimbra is a prime example.*)

A leading open source storage system platform is Nexenta (nexenta.org). A fully supported version is available from nexenta.com. This software runs on a Linux platform on a variety of commodity servers and supports a wide variety of commodity storage units. While the reliability and features of the product do not compare with the NetApp or EMC systems, it provides lower cost storage for some classes of storage.

Oracle's ZFS(open source)-based 7000 Storage Platform is a compromise between a fully open source DIY and a fully vendor supported "appliance".

Commodity Storage Devices

Full feature systems like NetApp or EMC, use higher cost disk components to meet their reliability life-cycle commitments. Such solutions cannot provide solutions that compete at costs the user community observes in the retail commodity market.

By giving somewhat on reliability expectations, and covering reliability with more backup, a storage service can be developed using "better quality" commodity storage.

An interesting compromise for higher density storage is the newer line of MAID (Massive Array of Idle Disks) which provides 80-120 TB of storage in a 6U shelf. Because the disks idle when not being used, power requirements are very low while access times are extended. (*BTW, powering disk on/off is not an issue ... consider laptop drives!*)

Cloud Storage

1. To meet the requirement for indeterminate-term archive storage, without increasing operating effort, CCS should consider "cloud"-based services. Cost effective for in-frequent access, these services would also include multi-site backup and media migration services. Contracts in the order of 5 years provides better pricing.
 - a. Initially at least, this "cloud" service should be implemented as a backend service to locally mounted services. e.g. A local Digital Archive on a MAID storage server could be backed up to the cloud service by CCS rather than signing many users up to the cloud.

Recommendations

Strategic Planning

Recommendation 1 that CCS develop a model plan for an Integrated Storage Infrastructure strategy based on a current industry model.

This plan must link institutional goals and policies with the technology roadmaps of the industry. Development of the plan will require some hard investigation of institutional requirements and

policy. Because most industry solutions are sector-focused and biased towards sets of requirements, a more complete understanding of vendor solutions and directions is required.

The combination of technical depth of technology direction and a strategic planning is not within CCS; the engagement of a strategic vendor can provide the capacity.

This is the #1 priority for CCS ... Making standalone tactical decisions to upgrade components of the storage environment may only complicate the ability to move to a comprehensive solution. Discussion of a strategic goal should be part of all upgrades and storage projects, rolling into a long term planning process.

Recommendation 2 that CCS encourage and support the development of institutional information management policies that clarify responsibilities for digital content management.

The organizational responsibilities for information stewardship is murky; Senior responsibility rests with the CIO but there is no process to reliability implement good stewardship policy in a distributed environment that increasingly relies on a CCS-provided storage services.

A starting point should be the development of an information classification model that permits discussion of requirements and permits users to consider if/how their information fits. A key part of the policy development then is to identify and classify the stewards that are responsible for policy development and implementation.

Tactical Direction

Recommendation 3 ... that CCS do a storage content analysis of enterprise and departmental storage to determine actual mix of storage classes.

Starting with the current content and available tools, CCS can develop a better understanding of the content now stored. Starting with the suggested classification in this document, the analysis can inform decision-making almost immediately. It is also important to consider regular audits that lead to a picture of the change happening in storage requirements as CCS clientele develops a higher dependence on an enterprise storage facility. The data collected will also be valuable in consultation with vendors to select and tune new systems.

Recommendation 4 ... that CCS immediately begin an RFI process that solicits recommendations on a suitable storage infrastructure plan and qualifies strategic vendors prior to an RFP for specific components.

As a starting point to building a strategic plan, the procurement of urgently required upgrades and more capacity should provide the opportunity to learn more about industry directions and identify potential strategic vendors.

A procurement project should start with a broader request for information from a number of leading vendors, asking not only for the storage components, but for a vendor description of how their recommendation fits with a recommended strategic direction. By providing some information on our data classification, the vendors can also identify ILM development strategies.

The result is an informed decision on a technology upgrade/expansion but also input to the development of a storage strategy.

Recommendation 5 ... that, as an interim step, CCS purchase a lower cost commodity SAN to extend CFS storage, implemented clearly as lower-performance, additional-drive storage with long but less frequent DR cycles. This solution would be replaced by the outcome of Rec. 4.

While the suggested data classification is “rough”, it is convincing that more low-cost capacity is required, and appropriate for the type of information being stored.

Since a lot of the content generation and storage process is still in the hands of individual users, leveraging CFS services is appropriate in the short term. Providing users with a secondary low-

cost, limited-performance, storage location while limiting growth of their primary store leaves the decision in the users hands ... for now.

Clear, pervasive, continuing marketing of the lower cost alternative is necessary for success but can provide a communications path for future enhancements and policy advice.

Recommendation 5 should guide the process of procuring this storage. Recommendations 6 and 7 describe parallel activities to expand the value of the low cost solution.

Recommendation 6 ... that CCS develop an information archive service, using existing storage capacity to begin to engage the community in information management and to off-load and re-direct storage requirements with lower access performance needs.

- A Research data store might be a priority special archive.

The Classification model suggests that there is a lot of content that users file away as a read-only object with a requirement of very infrequent access. Lifecycle is not considered. A web-based archive system which raises the questions about lifecycle and then uses the low cost storage system, provides an opportunity for users to think deliberately about lifecycle and to reduce the amount of “active” content they now try to manage,

From experiences of using internet-based archive services, once the user understands the concept, the continuing evolution of the archive system is not a problem as long as the information is readily accessible through a simple interface. CCS can begin with a rudimentary service application and evolve it or replace it without undue impact on users. Hence a “sooner-than-later” tactical plan is appropriate.

A big side benefit of an enterprise archive system is that the ILM process can provide valuable information on the characteristics (for classification) of the content. This project could provide early data for Recommendation #2.

Recommendation 7 ... that CCS initiate a project to investigate the opportunity/complexity of introducing lower cost open source and commodity components.

Examples include:

1. [NexentaStor](#) is an open source storage platform which can managed commodity storage units. A low-cost alternative to Network Appliance, it reduces the capital cost of a storage system by balancing off operating effort and reduced reliability (depending on configuration). Nexenta is used by a number of US Universities, primarily for lower cost, or second tier, storage. (e.g. see [Stanford story](#))

Both the Nexenta products and the Oracle ZFS-based storage servers were recommended as low cost storage alternatives by two Gartner consultants.

2. MAID (Massive Array of Idle Disks) is an emerging technology that couples requirements of low cost infrastructure with environmentally “Green” low operating cost. The technology is targeted at storage requirements of moderate-risk, limited-access, lower-performance, and high capacity. This appears to match the requirement of a large amount of infrequent-access and archived content. In conjunction with the development of a longer term strategy with a strategic vendor (Rec. #1), CCS could explore this technology for some specific applications.

In particular, storage of completed research data, as well as published content (e.g. institutional repository), would be a good fit. Backed up once for DR purposes, the MAID storage would provide great density and cost efficiency. When implemented behind a web-based archive or other application, the retrieval delays would be tolerated by users.

This project that would provide a CCS opportunity for exploration and innovation, encouraging the development of more depth of understanding storage needs and of the direction of storage technology.

Recommendation 8 ... that CCS implement a small scale “Permanent Historical Record” archive with a cloud provider to evaluate and demonstrate this opportunity.

There is a strong value proposition for the use of Cloud-based archive services, including avoidance of hassles with storage technology selection, rollover and infrastructure management. A full service provider would also provide the risk management for Business Continuity.

Choosing an archive service as an initial cloud storage application avoids the sometimes significant bandwidth costs to access the content. By its’ nature archived content is not accessed often. Data loss risks can be mitigated by continuing to maintain a local DR backup copy.

One of the issues arising in the community is that important historical information, now in digital form, is being lost because of neglect of the record. While stewardship policies and practices need to be developed, the Cloud Archive would be a good way for CCS to provision the infrastructure and some of the process while developing sound operating experience with a cloud provider.

Cloud services, as well as the technology, is evolving rapidly. Providers business plans must recover their costs as technology changes. Short term contracts attract higher charges; long contracts limit opportunity for competitive procurement. At the present time, it appears that a 5-year contract would attract the lowest costs while providing an opportunity to consider alternatives before content formats become stale.

Many cloud storage providers are targeting the individual consumer in the commodity market with smaller storage capacities and higher support overhead. For an enterprise, an alternative is a negotiated enterprise contract with a single point of contact. In this scenario, the University would provide the “front-end” archiving service locally. This server would stage content to the cloud service, reducing the provider complexity and support overhead. Access control then becomes a heavier requirement at the staging server.

An interesting combination is to implement a low cost local archive service based on an open source Nexenta platform using a MAID storage shelf and then “archiving” that content, with metadata, to a Cloud provider as a DR archive copy. CCS would be the only subscriber to that Cloud store.

However, if the service is constrained to an Institutional Historical Records application, the number of users involved would be fewer, and institutional curation procedures would integrate with access control.

Appendices

Appendix 1 Examples of lower-performance user storage challenges

1. **Accumulating old versions:** A departmental administrative office manages financial details in spreadsheets. Because more than one department manager needs access to the information, copies are distributed periodically by email. Each month, a snapshot of the spreadsheet is stored for backup and for month by month comparison. The unit has over 5 years of monthly snapshots since there are no clear policies on how long to maintain these records. The request is for low cost CFS shared space for the spreadsheet snapshots.
2. **Researcher Data – unfunded repository:** A researcher has completed several projects collecting data associated with a key research question. The researcher has also obtained related data from colleagues at other universities. This data, as well as intermediate processing results, are stored on their workstation hard drives. To analyze the data, components are copied to Sharcnet storage for HPC processing. Intermediate results are stored back on the workstation hard drive. The data is also analyzed by grad students who have copies on their hard drives. The researcher's main concern is that the volume is getting too large to manage and it is too time consuming to do backups. Research grants will pay for hard drives but not for long term storage of all the data. The researcher wants low cost (on department budget) enterprise storage for this read-only but growing data store.
3. **Collecting/preserving content for learning:** A faculty member provides a lot of content, including documents, data, and pictures, in the LMS for her students. Throughout the year she locates additional content related to the course, storing it on CFS and her workstation hard drive. Each semester, as the course interest evolves, she selects different content to load into the LMS. She also stores student submitted content, etc. from the LMS. The content volume is getting too large to manage and it is an effort to load and replace content in the LMS. The faculty member wants low cost CFS storage to house all her collected content or the ability to keep it all in the LMS.

Appendix 2 Rice University Storage systems Evaluation

This report is taken from the Rice University IT website. IT includes useful information about technology issues associated with selection of storage systems. In supporting a tiered storage plan, the inter-relationships between these different technologies need to be more deeply investigated as part of an RFI/RFP process.

Rice University IT Storage System Evaluations 2010/2011 Summary

Criteria

This evaluation was initially targeted at finding a viable and inexpensive solution to shared storage for virtual machines that could be easily managed, simple, supportable and have sufficient high availability to support core services that require high uptime.

History

Six years ago Rice enterprise storage consisted of an outdated and unsupported Auspex NAS that provided NFS mounts to a set of Unix file servers running Solaris, a number of Windows servers with direct attached storage shelves from the Compaq era and silos of SAN storage used to support legacy video archives and other services. We were able to consolidate and migrate almost all of these disparate systems into an enterprise SAN and Nas. The SAN consists of a mix of Hitachi and Enginio disks and the Nas is a mix of Enginio FC and Sata front ended by a Bluearc Nas gateway. We were looking for a highly resilient storage solution that was a dual system where data could be replicated or mirrored and fail over instantly as needed for maintenance (scheduled) or breakage (unscheduled) without disrupting services in the VM farm. Currently we have purchased under an evaluation a VM shared storage system that uses Nexenta ZFS as the core on SuperMicro storage components. The product evaluation had been a bit rocky to start, in that the system is not performing under NFS as expected. This problem was addressed by the vendor and it is providing good performance, ISCSI numbers are ridiculously high. We had found a number of issues all of which have been addressed at this time.

The investigations led us to evaluate a number of different solutions from inexpensive supported open source solutions to Cadillac vendors with high price tags and soup to nuts features. We evaluated both SAN, NAS, Hybrid and Direct Attached systems looking for a good fit. In the midst of this SSD drives were making their way onto the scene and causing a stir and threatening to completely disrupt most storage thinking if they can get a few things ironed out in the near future. This could be the equivalent of what disk drives did to magnetic tape and punch card technology as a storage media. There was also discussion regarding having one ring to rule them all – a single solution that could fulfill all of our storage needs. The prospects of having once piece of technology be the lynch pin that all services depend on was considered to be just too risky and no one wanted to hang their job on a single solution no matter how much the vendor tried to convince us that the second law of thermodynamics did not apply to their product. This does not mean we do not have faith in the vendors, we do, however we looked at what the cloud providers were doing and decided that a more modular approach was more prudent, inexpensive and future proof.

Types of systems

As you would expect, each vendor believed that the solution that they presented was uniquely qualified to provide a better service than their competitor, and we tried very hard to keep an open mind. We invited faculty, peers, administrators anyone who might be interested to help us discern what would be a good decision and fit for storage. In order to explain storage choices to the lay people in the groups that we talked to, it was necessary to segregate systems in a logical manner and here are some ways to break these down. The remainder of this document is devoted to explaining and classifying storage and breaking down the vendors we looked at.

Block based

SAN, iSCSI, FCOE and ATA/OE protocol based systems fall into this category and are systems that provide raw disk blocks to servers who can then make the storage available to end users or systems.

File based

NAS systems which traditionally deliver CIFS and NFS protocols fall into this category and are systems that provide file based services either directly to end users or to systems for use. To add confusion to the mix, most NAS systems can offer some block based protocols as well such as iSCSI and FCOE.

Chunklet based

These are Block based storage systems that don't use RAID to protect from disk failures. They spread multiple copies of blocks over the entire disk system so that you can have N layers of protection. The data is delivered from the chunklet blocks to an upper level system and distributed via block based or file based service. Some of these systems also had what is known as zero spares. This means that they don't recover from a disk failure by recovering the drive, they simply fail the bad drive and rebuild the data across all drives. The benefit being that it takes less than an hour to rebuild a 3 TB disk. Traditional spare drive systems can take many hours to days to rebuild a drive of this size while performance and safety are at risk during the rebuild.

I/O Performance mechanisms

Either Block or File base services have to perform to some minimum requirement based on the services that will be using them. There are a few tricks to getting performance and you have to understand the performance requirements of the consumers of your storage system in order to be able to provide the appropriate hardware.

1. **Read Performance:** If the consumers that will be using the storage service are read intensive, then caching systems may provide a good solution. These systems leverage inexpensive disks with low performance characteristics with SSD, Flash, battery backed NV Ram or other fast caching service in front of them. This can reduce costs by limiting the number of spindles required for storing data while providing good read performance. Unfortunately there are plateaus for write performance that may limit their use and both read and write performance needs to be considered.
2. **Write Performance:** Caching to flash or SSD drive technologies can supplement performance for writing to disk, but this is a stop gap. If your consumer applications do a large amount of writing data such as busy databases, email, etc. then you will need to be able to keep up or the applications will be slow. In order to do this traditionally you add a lot of disk drives (spindles) to your storage system and distribute the writes across many disks (AKA Striping DATA). In all of the improvements to disk storage technology, there has been a huge increase in the capacity of disk drives, but very little change in their ability to write more data on the disk faster which is only possible currently by spinning the disk faster and we have reached a technological barrier on that point. Write cache buffers on controllers, using flash or SSD can be used to head fake the applications and keep them from stalling as they wait for data to be written to the disk, but this only works if the amount of data out to the disks exceeds the amount of data into the cache. Once you over run the cache, your applications will stall. Again, you have to know what the

requirements of your applications are that are using the storage system in order to purchase the right system.

3. **Virtual Machines:** This technology has provided a wonderful ability to IT allowing us to segregate hardware from applications. As with any form of indirection, it also carries a measure of complexity. In our case, we have to keep watch on the number and type of servers guests that we load up on our VM hosts. If we don't, we could have problems with performance on the applications running on those hosts if we over subscribe the storage performance. Most of the storage services that we reviewed were approved by VMWare and some had VMWare triggering applications built in that would allow for dynamic provisioning and failover. Approval by VMWare did NOT include any performance characteristics.
4. **Storage Features:** As the price for disk drives have plummeted per Gigabyte of raw capacity, the cost of features has fought to hold storage prices high. Some of the features are in high demand for good reasons and you should pay close attention to these.
 1. **Thin Provisioning:** The ability to only allocate as much storage to a system or service as it is used or written on as opposed to pre-allocating a fixed amount that is mostly empty. This feature is available in both block and file based systems
 2. **High Availability:** The best type of high availability occurs when ANY system component or ANY service component fails and the service survives. Depending on your SLA, it may do little good to have a storage system failure bring down the services that run on it. Virtual machines aggregate many services into choke points. When these fail, many services go offline. HA can cause the costs of a storage system to go up dramatically – expect to pay more than double.
 3. **Primary De-duplication and compression:** These technologies have been around a long time, and are making a resurgence. Initially they targeted backups but have moved into cloud storage and primary storage. The focus is on correctly identifying identical blocks of data and only storing them once, then having the ability to recall and reconstitute the blocks in a meaningful way at the performance level required by the application. The goal is to reduce the amount of total storage space. Well there has to be a balance, you need spindles for write performance, but you want to reduce them to some sweet spot to keep costs low and disks are pretty inexpensive. So the question is do I spend a lot of money for a feature or do the cost of cheap disks override the complexity of de-duplication and compression. It should be noted that for backups, these technologies and encryption increase recovery times and you have to reconstitute any data to read it even if it is just for backups.
 4. **Centralized intelligence or distributed intelligence:** A centralized intelligence storage system may contain an expensive high throughput controller or head that everything flows through and can provide a lot of performance because it does most of its computing in custom chip sets that offload a lot of the work. Another model is to distribute the workload among many smaller compute engines (generally COTS PCs) that are responsible for segments of disks which are modular. Each of these has its merits and challenges in costs, features and scalability. The cloud model is the distributed intelligence (Google) and legacy enterprise is centralized intelligence.
 5. **Dumb storage and intelligent storage:** Dumb storage is what we call off the shelf JBODs that are inexpensive to buy and maintain and require intelligence provided by an upper level controller or manager. Intelligent storage is similar to modular systems where the controller is built into or provided with the storage. Again, intelligent storage has more features and costs more.
 6. **Dynamic Tiering:** A intelligent storage feature that moves data at the block level to lower cost / lower performance media. This is done on the fly via policy based decisions and provides benefit by optimizing or right sizing the amount of high performance vs low

performance disk assets. A good example is old data that is rarely accessed, why keep it on expensive high performance disk? This is similar to what caching does, but on a permanent storage level not as a transient indirection layer to provide better performance.

7. **VM Aware:** Some of the storage systems can intelligently communicate with VMWare management that allows for automation for Auto fail over and auto scaling.

Cloud storage

Cloud storage providers are available, and can provide a potential cost savings due to economies of scale for some amount of storage. The challenges with cloud are determining what is safe to put in a cloud, how to build a contract that will protect me if I have to change providers so that my data is not held for ransom, building a contingency should the cloud provider be bought or go out of business and making sure that critical data is not completely in the hands of a sole solution. Next you have to evaluate the costs for the storage, the costs to transmit the data and see if it works for your budget. In my view, the cloud is beginning to be cost effective and appropriate for backups to eliminate tape systems, tape management and to meet off site backup requirements and to provide storage to servers in a cloud service, but not as a primary data storage facility for a campus or campus systems. As this matures just a bit more, we see the possibility of cloud as primary storage front for many types of data with local SSD caches or possibly local low cost disk systems which deduplicate data compress and replicate deltas into the cloud.

Vendor Classifications

Vendor	Storage Type	Features	Protocol
Compellent (Dell)	Block Based - Chunklet	Thin provisioning, HA, Chunklet, Intellegent storage, Modular. Virtualization, VM Aware	Fiber channel
DataCore	Block Based	Thin provisioning, HA, Virtualization	Fiber channel, iSCSI
IBM	Block Based	Traditional SAN with Virtualization, HA, Thin Provisioning	Fiber channel
Hitachi	Block Based	Traditional SAN with Virtualization, HA, Thin provisioning	Fiber channel
Dell (EMC)	Block Based	Traditional SAN with iSCSI,	Fiber channel, iSCSI
Falconstor	Block Based	Virtualization, Thin Provisioning, HA	Fiber channel, iSCSI
NetApp	File Based	Thin provisioning, HA, Deduplication, Central intellegent storage, virtualization, VM Aware	NFS, Cifs, iSCSI, FCOE
NimbleStorage	Block Based	Deduplication, Thin provisioning, HA,	Fiber Channel, iSCSI
NorStore-Nexenta (Solaris 10 on Supermicro using ZFS)	File Based	Thin provisioning, HA, Deduplication, Modular, VM Aware	NFS, Cifs, iSCSI, Fiber Channel
3-Par (HP)	Block Based	Thin provisioning, HA, Modular, VM Aware, virtualization, Chunklet Based, Tier migration	Fiber Channel, iSCSI
Isilon (EMC)	File Based (chunklet)	Thin provisioning, HA, Modular, VM Aware, virtualization, Chunklet Based, Dynamic Tiering	NFS, CIFS, iSCSI
BlueArc (Hitachi/LSI/DataDirect)	File Based	Thin provisioning, HA, Central Intelligent storage, virtualization, VM Aware	NFS, CIFS, iSCSI
Coraid	Block Based	HA, Central Intelligent storage, VM Aware	ATAOE
Coraid-Nexenta (Solaris 10 on Coraid using ZFS)	File Based	Thin provisioning, HA, Deduplication, Modular, VM Aware	NFS, Cifs, iSCSI, ATAOE