

# Backup Workgroup Report

July 2010

Gaylon DeGeer, Steve Fowler, Clif Johnson, Andrew Morgan, Curtis Onstott, Sean SanRomani, Todd Shechter, Alan Sprague

## Table of Contents

Working Group Charge .....	1
Executive Summary.....	2
Backup Strategies.....	2
Data loss tolerance .....	2
Retention period.....	2
Recovery time .....	3
Backup Purpose .....	3
Retention .....	3
Tape versus Disk.....	3
Consolidation and Scaling .....	4
Conclusions .....	5
Recommendations .....	6
Appendix A – Current Backup Systems.....	7
Network Engineering .....	7
Enterprise Computing Services .....	7
College of Business .....	7
College of Engineering .....	7
College of Forestry.....	8
College of Science .....	8
Crop and Soil.....	8
Appendix B – Consolidation considerations .....	8

## Working Group Charge

1. Are our current backup strategies based on full and incremental backups to tape/disk using 3rd party backup software still appropriate given the inexpensive disk arrays currently employed on campus and the likelihood of an exponential increase in storage due to cheap disk availability?
2. Is there a disk based backup strategy that will still provide a reasonable level of availability but at a much lower cost? If so what are the trade-offs compared to our current approach?
3. Assuming a disk based strategy does make sense, is there a scale at which such a system could be leveraged by several campus groups in a cost effective way? Is there a limit to the scale at which such a system is cost effective?

## Executive Summary

At first glance the charge questions seem simple enough, but as the members of the working group dug into the question of backups, what we were doing in our individual areas, and the idea of a shared infrastructure it became apparent that there were a great deal of things to consider.

Backups are a necessary part of any IT infrastructure and every group has a backup process for their servers and data stores. The backup systems in use are all disk, all tape or a mix of tape and disk. The costs and capabilities of those systems vary depending on the vendor and features included in the system.

During our discussions we explored what each member group was doing for backup today including:

- Backup hardware and software used, including costs
- Backup retention periods and backup schedules
- The types and volume of data being backed up
- Unit policies and service level agreements

A summary of system information is provided in Appendix A.

We were asked to consider backup strategies based on the idea that “disk is cheap” and to consider whether there was a system size that made sense. As a group we didn’t find that there were viable options to current backup strategies as each group has worked to optimize the strategies used given system type, purpose and unit policies. As for costs, the price of disks has come down and the volume of data a disk can hold has gone up, but then so has tape. In an apples-to-apple comparison tape still edges out disk when cost is the only comparator. As to the question of scale, it is difficult to say where the cost-effective tipping point is due in part to the variety of hardware and software in use.

We don’t recommend a single monolithic backup system due to the potential costs and the logistics of managing such a system. We agree that some consolidation through inter-departmental or inter-college collaboration has merit and should be explored.

## Backup Strategies

When we looked at the issue of backup strategies we found that each group employs different strategies based on what they are backing up. Drivers or considerations for the different strategies are as follows:

### Data loss tolerance

How often does the data change and how often must it be backed up? Most backups at OSU occur daily except for a few archival file shares that are backed up quarterly. Some units employ shadow copies or snapshots (a feature of recent Windows systems and some UNIX file systems) to allow a point-in-time capture of data at even more frequent intervals, such as twice a day.

### Retention period

How long ago must data be available for recovery? Retention periods at OSU range from 1 month to 6 months. The retention period may be determined by regulatory requirements or service level agreements. See the Retention section for an in-depth discussion.

## Recovery time

How quickly can data be restored? Different technologies have different recovery times. Shadow copies and snapshots are the fastest to restore. Disk-based backups are slower but still quite fast. Tape backups are the slowest to restore.

## Backup Purpose

Are backups used only to recover files that are accidentally deleted or corrupted, or is disaster recovery required? For disaster recovery purposes, backups are typically stored in a separate building or off-site entirely (ship tapes to a vendor such as Iron Mountain).

## Retention

Department email, financial records, course grade lists, research data, and other types of data are only useful when they are secure from unwanted modification, corruption, or deletion, and can be accessed when needed. A retention policy describes the amount of time that data is required to be available to users who need access to it. A retention policy may also define when data should be destroyed. Generally, any data kept on network storage systems will be available for an indefinite amount of time. As storage technologies change, data is simply moved to the next technology and remains available. For most user data, this can easily be years of “retention” without interruption.

In backup software, the retention period is how long data is stored and therefore how far back in time a file can be recovered. The retention period should be long enough to prevent the loss of data that is modified, corrupted, or deleted. Shorter retention periods keep the costs of backups down. Longer retention times increase the burden on the backup system.

Typically, data is backed up according to a schedule of 1 full backup per month and 1 incremental backup per day. A full backup is a copy of the entire data set. An incremental backup is a copy of only the data that has changed since the last backup. As the retention period is increased, additional full and incremental backups must be stored. For a typical 6 month retention period, the backup server must store 7 full backups and 7\*30 days of incremental backups.

For example, let’s assume a data set of 200GB with 1GB of changes per day (assuming an average of 30 days per month):

Retention	Size of all Full backups	Size of all Incremental backups	Total backup storage required
1 month	400GB	60GB	460GB
2 month	600GB	90GB	690GB
6 month	1400GB	210GB	1610GB

**\*Note** – In order to recover data for the full retention period, an additional month of backups must be retained, which is why the sizes above appear to be 1 month too large.

As you can see, the retention period has a large effect on the amount of backup storage required. A moderate data set of 200GB is multiplied many times when the retention period is long.

## Tape versus Disk

Both technologies are currently in use at OSU in both small and large backup systems.

While we weren't asked to consider which technology was cheaper we think it is important to dispel the misconception that disk is significantly cheaper than tape.

Breaking down the component parts, if you were to search for prices of a 1 terabyte disk and an LTO 4 tape or a 2 terabyte disk and a LTO5 tape you'd find that they are fairly close in price. Likewise if you were to price enclosures for the two technologies you would find a wide range of prices, both about the same depending on the vendor and features.

Product *	Capacity (Unformatted)	Cost	Cost/GB
Hard drive SATA	1TB	\$80 - \$200	\$.08 - \$.20
Tape LTO4	800GB/1.6TB**	\$31 - \$94	\$.031 - \$.094
Hard drive SATA	2TB	\$130 - \$370	\$.065 - \$.185
Tape LTO5	1.5TB/3.0TB**	\$105 - \$126	\$.052 - \$.063

\* All media are consumer grade products with prices reviewed on Nextag.com. Enterprise class disks such as Serial Attached SCSI (SAS) or Fiber Channel (FC) disks jump to \$.70 - \$1.00/GB in cost.

\*\* Capacities reflect uncompressed and compressed data capacities.

Then why is disk being used in place of tape?

Disk Advantages:

- Faster restores due to random read capability of disk and data is online at all times
- No need for an operator to swap tapes
- Concurrent read/write
- Faster writes for large numbers of small files
- Able to handle a greater number of backup clients simultaneously and meet backup window demands.
- TCO Advantages
  - Less time to manage – no tape swapping and cataloging
  - Lower licensing costs compared to large tape libraries (Specifically EMC Legato)

Tape Advantages:

- Easier to store backups off-site by moving tapes
- Highly scalable - additional tapes can be purchased to expand capacity with very little cost
- Lower cost for multiple redundant backup sets
- Faster writes for large files
- If a tape is corrupted or damaged the same data may be on multiple tapes. If a disk array is damaged then all data may be lost.
- TCO Advantages
  - Lower power/cooling requirements thus a "greener" alternative to disk.

## Consolidation and Scaling

We investigated the feasibility of consolidating several backup systems in order to save costs. As the number of clients and the size of backups increase, the complexity and cost of the consolidated backup

system actually increases. Some specific issues that would have to be considered in a monolithic system include network throughput, backup window management, number of systems that can be serviced in a given backup window, the number of backup schedules/retention intervals that can be offered, system resiliency, total data storage needed, and system management responsibility. We also considered the implications of a single system if the unthinkable happened and everyone was dependent on this one system to restore their infrastructure.

The minimum cost to build a backup system is approximately \$25,000. This includes a backup server, backup software, and a disk or tape storage system. One of these “starter” backup systems can probably handle up to 10 clients and store 20-30TB of backup data (don’t forget the retention multiplying effect). For disk backup, 20-30TB can be consumed very quickly especially if data is retained for several months.

Additional clients or larger data sets would require a larger disk storage system or additional tape libraries and tape drives. For instance doubling the number of clients or the backup size will require twice as many tape drives, tape libraries, and/or disk storage. As the system scales up a certain point is reached where additional backup servers are required to handle the aggregate throughput of all the backup clients. The other thing that increases is the licensing fees for the backup software. Some of this can be off-set when using disk as a backup medium by employing data de-duplication technology. The trade-off of using de-duplication is that the acquisition cost goes up significantly.

## Conclusions

*1. Are our current backup strategies based on full and incremental backups to tape/disk using 3rd party backup software still appropriate given the inexpensive disk arrays currently employed on campus and the likelihood of an exponential increase in storage due to cheap disk availability?*

Each backup operator has evaluated their current backup strategy to determine if it is meeting their requirements (see the Backup Strategies section). A full and incremental strategy is not used in every case, but only where the data loss tolerance demands its use. In some cases shadow copies or snap shots are used to further enhance and shorten the data loss window. Yes, the full and incremental strategy is still appropriate and we would go so far as to say it is absolutely essential in any disaster recovery plan where no more than a day’s worth of data can be lost.

Using an inexpensive disk array is not always appropriate. This is especially true where we have multi-month backup retention. In that situation inexpensive disk arrays have serious short comings. For example, 2TB of shared files with monthly full backups, daily incremental backups, and six month retention can consume more than 40TB of disk space. This is dependent upon the percentage of daily data change and won’t be true in all cases, but this particular example comes from actual usage figures of systems here on campus. It should be apparent that a large number of inexpensive disk arrays would be required given the amount of data currently being stored without some additional technology such as de-duplication.

*2. Is there a disk based backup strategy that will still provide a reasonable level of availability but at a much lower cost? If so what are the trade-offs compared to our current approach?* The group wasn’t sure what was meant by “level of availability” and was provided the following by Cheri Pancake: *My understanding of the question about "availability" was that it included both frequency of backup and*

*retention time - and whether one-level-fits-all or whether there need to be different tiers for "critical institutional information" as opposed to "run-of-the-mill files."*

We think it is safe to say that there isn't a one-size-fits-all solution unless everyone was to agree to the same set of constraints. From the technician's point of view one month retention is perfectly reasonable given the number of systems and quantity of data that needs to be backed up. From the customer's point of view this might seem unreasonable especially where expectations of a longer retention period have been in play. It is obvious that the shorter the retention period the smaller the backup system has to be and thus the lower the cost of the system overall.

Cherri brings to light another issue and that is the difficulty of determining data "value". It is up to each user to make this value judgment. To accommodate users some groups provide users with data space that isn't backed up at all. This approach helps to lower the cost of storage by removing the cost of backup, but it also carries some element of risk because if something catastrophic happens to that storage then all data is lost.

The trade off to all of this would be that users could no longer expect to have files recovered up to six months after they were deleted.

*3. Assuming a disk based strategy does make sense, is there a scale at which such a system could be leveraged by several campus groups in a cost effective way? Is there a limit to the scale at which such a system is cost effective?*

Without designing a system and making some assumptions about the numerous variables that have to be considered in such a system we are not able to speak to the question of scale. Using Network Engineering's system as a possible example, not counting FTE costs, this disk/tape hybrid system cost close to \$200k and has a recurring annual cost in excess of \$30k in maintenance/support to software and hardware vendors. In April it backed up 42TB from 100 client servers owned by 16 different departments. Any major addition of either clients or data would require significant capital investments in hardware and software to accommodate that addition.

We don't believe a disk based strategy make sense. Several of the participating units are retaining their tape systems because tape is performing adequately and migrating from tape to disk would require a large initial investment to support their current backup strategy.

## Recommendations

We do not recommend that OSU move to a single monolithic backup system. The cost of designing and implementing such a system could exceed what groups have spent on their existing systems. It also creates a single choke point if multiple users are trying to restore systems at the same time.

Consolidation of smaller, individual backup groups does have merit. It spreads the cost of the system out over several departments, eliminates the single choke point, and spreads the clients out which minimizes contention and scheduling issues on the backup systems. Having several backup centers might also allow sharing of resources. One thought was the ability to share reserve capacity in an emergency. If a situation arose where the primary backup system was not available, some critical systems could be backed up to another center.

We also recommend that OSU look at trying to get better pricing for Networker Licensing. This product is used in the majority of large backup systems thus representing considerable cost to the university. This cost is also one reason smaller units have opted to go with open source software or a different backup software vendor.

## **Appendix A – Current Backup Systems**

### **Network Engineering**

Network Engineering runs a centralized backup system for several departments, although it is mainly used by central IS. The system is funded by per-GB fees for backups and located in Dearborn. The backup server is a Sun v440 running EMC Networker backup software. There are 2 tape libraries and a disk storage system connected to it. The Spectra64 AIT4 tape library is being phased out due to its high annual maintenance costs and older AIT4 tape technology. A new Overland NEO 2000 tape library, which contains 30 tape slots and 2 LTO4 tape drives, was purchased this year. There is also a Data Domain DD580 with 15TB of raw disk storage. The DD580 employs data de-duplication software, which gives the system a logical capacity of 150TB based on our current backups.

### **Enterprise Computing Services**

Enterprise Computing Services runs a remote storage node that is connected to Network Engineering's EMC Networker backup system. The storage node (server) is an HP DL380G5. Attached to it is an HP MSL4048 tape library which contains 48 tape slots and 2 LTO4 tape drives. This system is strictly for Banner backups. Tapes are sent off-site to Iron Mountain.

### **College of Business**

The College of business currently utilizes Symantec NetBackup Enterprise as our primary backup service. We are running it on a single Dell R300 which is both the master and media server. The backup library is a Dell TL4000 robotic library which has 48 slots and 2 of the available half height drive bays populated with LTO 3 drives. We are in the process of adding an older EMC CX300 for some limited Disk-to-Disk backups and replicated storage for our primary file server.

### **College of Engineering**

The College of Engineering uses EMC/Networker software running on a Solaris server as our main backup host. Attached to the Solaris server via 10G fiber are two DataDomain de-duplication appliances, each with their own set of hard disks. Engineering also has a large tape library that it uses for archiving data, longer term storage. Engineering has a 3 month retention policy for current data on the DataDomain equipment. The funds necessary to run Engineering's backup system come from both Engineering IT's general operating budget and from researcher contributions.

## College of Forestry

The College of Forestry uses EMC/NetWorker software to manage the backup system. The system is built on a Dell PE2900 with a Qualstar TLS-412180 AIT-3 Library, with 10 tape drives. In addition there are two storage nodes in the system, one with a Qualstar TLS-46120 AIT-4 Library with 6 tapes drives, and one with a Qualstar RLS-4221 AIT-5 Library with one tape drive. General costs are built into user and group fees and tired based on backup strategy and retention period requirements. Additional backup and retention strategies are available and costs are based on retention periods and total GB backed up. With AIT being at its end of life, replacement hardware solutions are being investigated.

## College of Science

The College of Science has a centralized backup system for all COS departments, and paid for research servers. The system is an HP DL380 server, running Linux and the open source backup software Bacula ([www.bacula.org](http://www.bacula.org)). The backup system is connected via NFS to a Silicon Mechanics 24TB NAS, which provides the COS with a 30 day retention period for disk-to-disk backups/restores. The COS also owns a NEO 2000 tape library, with 2 SDLT tapes, and 26 slots. This unit is currently in need of repair, but could be utilized to extend backup retention periods, increased backup capacity and offsite storage. The backup server, NAS and Tape library are located in a building on campus separate from where the servers which are being backed up are located.

## Crop and Soil

Crop and Soil Science uses Symantec Backup Exec running on an IBM X3400 attached to a Qualstar 4222i 20 slot changer with an AIT-5 drive. It also has a dual port Qlogic ISCSI HBA attached to an AC&NC external Raid using 2 1gig links with 16 raw Terabytes of storage. We run it at RAID level 6 with a hot spare so actual capacity for b2d is about 12 terabytes. We will soon be swapping this out for an identical AC&NC Raid with 32 raw terabytes. We do a mixture of b2d and tape. Tapes are stored in a Safe on-site with a 3 month retention period. Costs for this system have been shared with the CAS deans office since 1996 and startup costs are in line with the estimates above.

## Appendix B – Consolidation considerations

Here are several factors that should be considered when working toward consolidating backup systems. These are in no particular order as some factors may be more or less important depending on the systems being consolidated.

1. Total backup volume. This is the amount of data backed up times the period the data is retained plus any projected growth.
2. Number of systems being backed up. There is a limit to the number of systems that can be backed up in a given period of time. This is usually determined by the number of backup threads the backup server can handle, and the amount of data to be moved and indexed.
3. Backup strategies. Not every system needs to have full plus incremental or differential backups. Some systems may only require a single monthly full.

4. Network. The speed of the network connection can be a choke point if the speed at which data is moved isn't fast enough to allow the backup to complete before the next backup window starts. The entire end-to-end network path needs to be considered not just the capability at the server. Firewalls can also be an issue if the proper ports aren't allowed.

5. Cost sharing.

6. System management/responsibility. Some control has to be relinquished and a level of trust placed in the organization that is hosting backup system.

7. System security and location. The space where the backup system is housed should be a secure, controlled access room. It is preferable to have it in a location other than the same room with the backed up systems.

8. Backup software. How much administration can be delegated, which operating systems are supported, and how does it handle open files and which ones are affected?