



Global Information Assurance Certification Paper

Copyright SANS Institute
Author Retains Full Rights

This paper is taken from the GIAC directory of certified professionals. Reposting is not permitted without express written permission.

Interested in learning more?

Check out the list of upcoming events offering
"Security Essentials: Network, Endpoint, and Cloud (Security 401)"
at <http://www.giac.org/registration/gsec>

Mike Taylor
V1.2e

Centralized network backup of UNIX systems to a SpectraLogic AIT-2 tape library utilizing the AMANDA archiving utility.

Introduction:

Confidentiality, Integrity and Availability is the mantra all of us involved with data/computer systems security are chanting. Proper backups of data can help to ensure both the integrity and availability of data stored on your systems.

AMANDA (Advanced Maryland Automatic Network Disk Archiver) is a backup system designed to back up multiple servers or workstations to a centralized backup server with a large capacity tape drive. The Spectra Logic company designs and manufactures automated tape libraries and other products that protect customer data. AIT-2 tape technology is Sony's latest 8mm product delivering transfer rates up to 12MB/sec and capacities up to 100GB on a single tape. The Spectra Logic Bullfrog system, which can support up to four AIT drives and 40 tapes per unit, is the system that will be discussed here. Spectra Logic does not directly support the use of AMANDA software. Anyone wishing to utilize AMANDA with a Spectra Logic Bullfrog Library will be "on their own". The remainder of this document will attempt to explain some of the special issues you may need to address when considering such a solution.

Why use a Tape Library:

The Requirements for network storage have increased dramatically over the past years and there is no indication that this trend will change. While the requirements for data storage have increased, the resources available to address these issues have not increased at the same rate. The first method used to address this issue was to simply increase the capacity of individual tapes, allowing the backup of an entire network in one evening, without requiring an operator on duty to change tapes. However, this was just a temporary solution as most modern networks and some servers cannot be backed up to a single high capacity tape. Some type of backup automation now becomes necessary.

Backup automation serves to increase the backup capacity available without operator intervention. A sequential access device, commonly referred to as a stacker, can meet this need at a lower cost than a random access library. A stacker commonly features a single tape drive and just enough robotics to load the next tape. Usually there is no way to select a specific tape or to backup to a previous tape without manual intervention.

While the justification for backup automation is usually increased capacity, the justification for a Random Access Library can be to eliminate the human factor from the backup equation.

A properly sized tape library should be able to provide the required backup services with little or no operator intervention. Not only does this decrease the labor costs associated with routine tape changing tasks, it also reduces the chance for operator error. Once properly configured, the tape library will not use the wrong tape for a particular job or fail because the right tape was not in the drive at the time the backup (or restore) was requested.

Another benefit to utilizing a tape library for backups is centralized control of data. This can allow a system administrator, by aggregating the data and media, to more effectively enforce backup policies and procedures.

Also tape storage and retrieval can often introduce undesirable delays in the recovery of data, the ability to access the most recent data without having to retrieve the tape from off-site storage can be a great advantage when the dealing with mission critical data.

Why use AMANDA, instead of one of the Commercial backup packages:

Short answer, price. Most commercial backup software packages have a relatively high initial purchase price and often times have restrictions on the number of nodes that can be backed up. Adding additional nodes can be both expensive and time consuming, due to the need to secure additional funding and the associated delays in license updates. Furthermore, commercial packages are usually distributed in binary form for a specific operating system and version, which can add additional costs and problems when dealing with multiple operating systems and differing versions of the same operating system. Add to that the additional costs for drivers for your library that have been certified by your software vendor, and the price can be as much or more than you initially invested in your tape library.

An example of this pricing dilemma is illustrated here. Solstice backup (AKA Legato Networker) lists the price for the Autochanger module for 1-64 slots at \$8350.00 USD and the listed price for 100-node license is \$18000.

AMANDA however, is provided for free to anyone who wishes to attempt to use it. It is provided in source form and can be compiled on multiple architectures without restriction. However, AMANDA is provided “as-is” and as a result there is no formal support other than the “Amanda-users” mailing list, this can be an issue in shops where the internal expertise is limited.

Great I'm sold, what do I do now:

One of the first things you should do is review the articles at the SANS website. At the time of this writing there are two articles at that site that contain good information and links to other sources of information on this program. The articles are titled “Amanda, the Advanced Maryland Automated Network Disk Archiver. By Drew Einhom” and “Using AMANDA for High Performance Backups. By Laurence G. Guentert.” While I admittedly did not use these references in my initial setup of AMANDA, reviewing them could save you valuable microseconds during your initial configuration.

Purchase a Tape Library:

The selection and purchase of a tape library can be a little intimidating for the first time buyer. These units are not cheap; pricing starts at around \$20000 USD for the entry-level libraries, and this does not include the tapes. AIT-2 tapes, at the time of this writing, are selling for between \$85 and \$110 each. It is entirely possible that your future at your place of employment could be adversely affected by incorrect decisions made at this stage of the process. At the very least, making the wrong decision could be embarrassing and could affect your future credibility.

Spectra logic tape libraries have 2 different types of connection options, Direct Attach SCSI and Fiber Channel. SCSI is offered in both HVD (High Voltage Differential) and LVD (Low voltage Differential). The choice you make here will be largely dependant on the hardware available for the server you are using as a "Backup Server". For my purposes I used a Sun E3500 server, which had 2 separate fast-wide SCSI interface connections available extemally. This decision to use an E3500 (an admittedly oversized server for this purpose) was purely based on convenience, as any other server in my collection would have required additional hardware. It should be noted at this point that the SpectraLogic library requires 2 SCSI connections in the 3 or 4 drive configuration, however in the 1 or 2 drive configurations only 1 SCSI connection is required.

The number of drives required in your library configuration is a purely mathematical function (or should be, financial considerations may overrule your desire to backup everything in less than an hour). To determine the required number of drives in your library start with the total amount of data to be backed up for all servers and divide that by the total time you have available to perform the backup. This will give you the required data transfer rate. For example lets say you have 500GB of data you need to back up in 5 hours or less. This gives you a required data transfer rate of 100GB/hr or 1666Mb/minute or 27.8MB/second. While an AIT-2 drive is rated at up to 12MB/Second, my experience is that the average transfer rate is about 8MB/second. At 8MB/second you would need 3.5 drives, which we would round up to 4 drives. So, in this example we could, in theory, dump 500GB of data to 4 drives in less than 5 hours. That's just theory though. In practice you will very rarely get the full 100MB per tape, and changing tapes takes a few minutes, but this general guideline will be good enough for most installations.

The number of tapes required is a little more complicated. In the above example the pure math says 500GB only requires 5 tapes. But that's assuming that you will get the full 2:1 compression ratio. In reality you will find that most of your data is already compressed, and adding compression to already compressed files tends to make them LARGER, and compression tends to slow things down a bit. For capacity planning purposes use the "native" size. For AIT-2 tapes the native size is 50GB, so a full dump of 500GB would require 10 tapes. Now things start to get interesting (or silly). To do a full dump of 500GB every day, and have them available in the library for 30 days would require a library of 300 tapes. At roughly \$100/tape you would be looking at \$30000 just for tapes, and that would not cover the tapes required for off-site storage. Even if you only moved a full dump off site once a month, that is 120 tapes per year, in our example, that would be required for off-site storage. That's just way too many tapes, the chances of something happening to a set of tapes without being noticed increases with the number of tapes, so we want to keep the number down to a reasonable level to decrease the possibility that a set could "walk-off" without being noticed, a severe breach in confidentiality. Also, the Spectra Logic library that we are looking at only holds 40 tapes, not 300.

So, what is the solution to this dilemma (no, its not to make dilemmanade). The answer is "incremental backups". Incremental backups are a procedure in which after a full backup is preformed only the files that have changed since the last backup are scheduled for the next backup. In most cases the data that has changed is only a small percentage of the total data. Consider the following example. In this excerpt from an AMANDA report, the size of a full backup for a particular partition is listed.

DUMPER STATS

TAPER STATS

HOSTNAME	DISK	L	ORIG-KB	OUT-KB	COMP%	MMM:SS	KB/s	MMM:SS	KB/s
sapphire	/www	0	1699071	359264	21.1	26:50	223.2	1:12	4956.7

Comparing this to the output for an incremental dump of the same partition there is a noticeable difference.

		DUMPER STATS			TAPER STATS				
HOSTNAME	DISK	L	ORIG-KB	OUT-KB	COMP%	MMM:SS	KB/s	MMM:SS	KB/s
sapphire	/www	1	1026303	138784	13.5	12:43	181.9	0:30	4588.2

As you can see (by punching the numbers into your calculator) that almost 700MB had not changed and therefore was not backed up. While this may not seem like a big difference when a single tape can hold 50GB, but this was on a fairly active partition. Consider the partitions that never, or very rarely change. The following example illustrates this point. In the first report excerpt you can see the size of the full dump is listed as almost 2GB.

		DUMPER STATS			TAPER STATS				
HOSTNAME	DISK	L	ORIG-KB	OUT-KB	COMP%	MMM:SS	KB/s	MMM:SS	KB/s
zirconia	/usr	0	2458815	740640	30.1	22:31	548.1	2:11	5660.3

However during the incremental backup the dump size dropped to a mere 66MB.

		DUMPER STATS			TAPER STATS				
HOSTNAME	DISK	L	ORIG-KB	OUT-KB	COMP%	MMM:SS	KB/s	MMM:SS	KB/s
zirconia	/usr	1	66847	7488	11.2	0:34	221.4	0:03	2919.2

Now you are thinking “this is all great and wonderful, but 4 drives, 10 tapes, 5 hours, this is still going to get ugly when I try to do that full backup”, and normally, with most backup software this would be true. Now we can talk a little bit about the magic of AMANDA.

AMANDA Scheduling and Multiple Tape drives.

AMANDA is unique in its scheduling methodology. Unlike other systems which have very detailed scheduling abilities allowing you to specify when to perform incremental backups and when to perform full backups, AMANDA does not do this. Instead AMANDA only wants to know 3 things. How many tapes are available, how many days are in a backup cycle, and how many full dumps of each partition do you want during a backup cycle. This is a little confusing at first, until you understand that AMANDA is meant to free you from the trivialities of exactly when incremental or full backups are performed and will attempt to optimize each tape dump. Below is an excerpt from an AMANDA configuration file.

dumpcycle 1 weeks # the number of days in the normal dump cycle
 runspercycle 7 # the number of amdump runs in dumpcycle days
 tapecycle 40 tapes # the number of tapes in rotation

What this says is a full dump is performed a minimum of 1 time every 7 days for any given partition and 40 tapes are available for this process. The tricky part is that the full dumps are not all performed on the same day; they are spread out over the seven days to more evenly utilize the tape capacity. This may be a concern to some, but in reality which do you care most about, making the most of you tape backup system, or scheduling partition X for full backups on day Z. I have found that by letting the software manage the backups that the tapes are fairly evenly utilized. Here are the statistics from two different days to illustrate this point.

STATISTICS:

	Total	Full	Daily	
	-----	-----	-----	
Tape Time (hrs:min)	0:25	0:16	0:09	
Tape Size (meg)	10221.5	5433.2	4788.3	
Tape Used (%)	18.8	10.0	8.8	(level:#disks ...)
Filesystems Taped	25	9	16	(1:15 2:1)

STATISTICS:

	Total	Full	Daily	
	-----	-----	-----	
Tape Time (hrs:min)	0:15	0:08	0:07	
Tape Size (meg)	6645.4	3627.5	3017.9	
Tape Used (%)	12.3	6.7	5.6	(level:#disks ...)
Filesystems Taped	25	13	12	(1:10 2:1 3:1)

As you can see from this example although the number of file systems dumped to tape on the two separate days did not change, the number of full backups and the number of incremental backups are different, however the total percentage of tape used is very close between the days, showing the “averaging effect” that occurs when using AMANDA. This type of averaging that AMANDA performs automatically would be very difficult if not impossible to accomplish using commercial software solutions.

In these examples we are focused on a single tape drive. AMANDA does not understand multiple drives, however, it does give you the ability to specify a configuration file on the command line. Using multiple configuration files with small modifications can allow you to configure AMANDA to utilize multiple drives. Here is an example of a configuration directory structure setup for multiple drives, in this case two.

```

flash:/usr/local/etc/amanda#
drwxr-xr-x 5 root  other   512 Jun 26 16:56 ./
drwxr-xr-x 3 root  other   512 Nov  1 2000 ../
drwxr-xr-x 5 amanda other  3072 Jun 26 01:54 DailySet1/
  
```

```
drwxr-xr-x 5 amanda other 2560 Jun 26 02:33 DailySet2/
-rwxrwxrwx 1 root other 3 Jun 26 01:20 tape0-slot*
-rwxrwxrwx 1 root other 3 Jun 26 01:30 tape1-slot*
```

The two configuration directories DailySet1 and DailySet2 contain identical files, with the exception of the disklist file, which lists the servers and partitions to be backed up, and the chg-scsi.conf file which specifies which tape drive to use for this configuration. This is the part where a little advanced planning can go along way. To get the most out of your backup system try to split the partitions evenly between the two (or three or four) disklist files to help evenly spread the load between drives. While AMANDA easily balances the load between tapes, multiple instances of AMANDA are unaware of each other so balancing the load between drives is an exercise for the user. A little trial and error is necessary here as servers tend to vary in their response times and you are still trying to get all of them backed up in a small backup window. One thing worth noting at this point is that multiple instances of AMANDA cannot be started simultaneously they will error out while fighting over control of the robotics. Here is the script I use to run start amanda, which is run nightly from cron.

```
#!/bin/bash
#
# Empty the drives .
#
cd /usr/local/etc/amanda/DailySet1
/usr/local/libexec/chg-scsi -eject
#
# Pause to let the robot settle
sleep 120
#
cd /usr/local/etc/amanda/DailySet2
/usr/local/libexec/chg-scsi -eject
#
# Check then start first backup set
#
/usr/local/sbin/amcheck -m DailySet1
#
/usr/local/sbin/amdump DailySet1 &
#
# Pause to let the robot settle
#
sleep 480
#
# Next we check then start second backup set
#
/usr/local/sbin/amcheck -m DailySet2
#
/usr/local/sbin/amdump DailySet2 &
```

#

In the first section of the script I eject the tapes from both drives and put them back in the rack. I have found that the robotics/AMANDA cannot always find the tape if you left it in the drive. Next I check the configuration and start the first backup set. Then a relatively lengthy pause is needed before starting the next backup set. The pause is necessary due to the fact we have multiple drives but only one robot to serve them, also, our pointer to which tape to use next is shared between all drives and is not incremented by AMANDA until after the backup starts. Failure to wait will result in each instance of AMANDA trying to find and load the same tape into each drive. This will cause a failure of your backups. Once each drive is loaded with the proper tape and the dumps have begun multiple instances of AMANDA will run smoothly, each generating a separate report upon completion.

Requisite summation:

Backups are an essential part of any security program. Backups help to ensure Availability and Integrity of data. Poorly managed or overly cumbersome backup systems lend themselves to the possibility of breaches of confidentiality. Tape libraries help mitigate the risk of breaches of confidentiality by removing, or reducing human contact with portable media. Properly chosen and configured hardware and software can turn a backup nightmare into something manageable by mere mortals.

© SANS Institute 2000 - 2002
All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, or by any information storage or retrieval system, without the prior written permission of SANS Institute.

REFERENCES:

“What Is AMANDA” <http://www.amanda.org/>

“Spectra Logic – A brief overview of our company.”
<http://www.spectralogic.com/company/index.cfm>

“Spectra Logic – Media Technologies”
<http://www.spectralogic.com/technology/media.cfm>

“Bullfrog Library Family Overview”
<http://www.spectralogic.com/products/bullfrog/index.cfm>

“Compatibility Matrix”
http://www.spectralogic.com/common/collateral/interop/ISV_BF.pdf

“Solstice Backup Autochanger Software Module Licenses”
<http://store.sun.com/catalog/doc/BrowsePage.jhtml?cid=56465>

“Solstice Backup Client Connections”
<http://store.sun.com/catalog/doc/BrowsePage.jhtml?cid=56466>

“Interface Options” <http://www.spectralogic.com/technology/interfaces.cfm>

Amanda, the Advanced Maryland Automated Network Disk Archiver
Drew Einhom
December 19, 2000 <http://www.sans.org/infosecFAQ/incident/amanda.htm>

Using AMANDA for High Performance Backups
Laurence G. Guentert
January 29, 2001 <http://www.sans.org/infosecFAQ/unix/amanda.htm>

“Using Tape Libraries”
http://www.spectralogic.com/common/collateral/whitepapers/White_Paper_Using_Tape_Libraries.pdf