



INFN/TC-00/10

1 Giugno 2000

Backup dei dati utente nell'INFN: requisiti e valutazione della tecnologia hardware e software

Giovanni Barbagallo¹, Massimo Carboni², Davide Cobai³, Roberto Ferrari⁴,
Francesco Ferrera¹, Michele Gambetti⁵, Francesco Prelz⁶,
Ivo Saccarola⁴, Claudio Strizzolo⁷, Lucio Strizzolo⁷

¹⁾ INFN, Laboratori Nazionali del Sud, ²⁾ INFN, Laboratori Nazionali di Frascati,
³⁾ INFN, Gruppo collegato di Udine, ⁴⁾ INFN, Sezione di Padova,
⁵⁾ INFN, Sezione di Ferrara ⁶⁾ INFN, Sezione di Milano ⁷⁾ INFN, Sezione di Trieste

Riassunto

In questo documento vengono descritti i requisiti per un sistema di backup comune per i servizi calcolo INFN identificati dal gruppo di lavoro "Tool di interesse generale" su mandato della Commissione Calcolo.

Vengono inoltre illustrati i risultati e le conclusioni ottenute nel corso delle attività di valutazione di pacchetti software commerciali e *public domain*. Sono infine descritti i risultati delle attività di *technology tracking* di supporti e soluzioni di storage su disco e nastro.

PACS:89.80

*Published by SIS-Pubblicazioni
Laboratori Nazionali di Frascati*

Capitolo 1

Introduzione. Processo di analisi del problema

La prima applicazione sulla quale, su richiesta della Commissione Calcolo INFN, si è impegnato il gruppo di lavoro sui “Tool di interesse generale” è il backup dei dati di utente nelle sezioni.

Il processo di analisi e soluzione del problema ha riguardato i seguenti punti:

- raccolta preliminare delle esigenze delle sedi INFN attraverso un sondaggio;
- analisi e stesura dei requisiti del sistema;
- raccolta dei commenti sui requisiti sulla *mailing list*: `calcolo@infn.it`;
- valutazione di prodotti commerciali e soluzioni *public domain*;
- attività di *technology tracking* sulle soluzioni ed i supporti disco e nastro;
- scelta di una soluzione per il servizio di backup;
- configurazione di un *test-stand* per la soluzione scelta e prova intensiva della soluzione;
- elaborazione di un documento di raccomandazione della soluzione hardware/software da adottare, con note di installazione e configurazione.

Capitolo 2

Raccolta preliminare delle esigenze delle sezioni

La prima attività svolta è stata la raccolta, attraverso un sondaggio rivolto agli amministratori dei servizi calcolo, dello stato delle esigenze delle sezioni e delle soluzioni attualmente adottate per il backup dei dati utente. Presentiamo un sommario dei risultati, premettendo che al sondaggio, svolto nel mese di Giugno 1999, hanno risposto solo 16 delle 31 sedi INFN.

In primo luogo si è valutata la “provenienza” dei dati di backup, cioè su quali piattaforme e filesystem risiedono i dati utente oggetto del backup:

Provenienza dei dati di backup	Percentuale delle sezioni che dichiarano di effettuare il backup dei dati con questa provenienza, sul totale di quelle che hanno risposto.
Filesystem montato su Windows NT	40 %
Filesystem montato su Windows 9X	33 %
Filesystem montato su UNIX/Posix	100 %
Filesystem locale di sistema MAC	13 %
Filesystem esportato con NFS/AFS	60 %
Filesystem esportato via Netbios	13 %
Filesystem esportato con Appletalk	0 %

Questo il volume *massimo* dei dati da trattare:

	Volume massimo di backup
giornaliero	15 GB
settimanale	180 GB
mensile	100 GB
occasionale	15 GB/settimana

In media, il 55 % di questi dati risiede sulla LAN, altrove rispetto al sistema di backup. In alcuni siti più del 90 % dei dati di backup è remoto.

È stato poi richiesto quale software commerciale viene eventualmente utilizzato per il backup: 4 sedi utilizzano il sistema Legato Networker, 1 sede utilizza HP Omniback.

L'ultima domanda riguardava la destinazione dei dati: si tratta sostanzialmente di unità singole o robot DLT e unità DAT. 12 sedi su 16 (75 %) sono già dotate di unità singole o robot DLT.

Capitolo 3

Requisiti della soluzione di backup

Nella riunione del gruppo di lavoro del 15 giugno 1999 sono stati definiti i requisiti della soluzione di backup, basati sui risultati del sondaggio di cui sopra, e che sono poi stati inviati alla lista `calcolo@inf.n.it` per raccogliere commenti e proposte. I requisiti sono stati suddivisi nei due livelli di priorità seguenti:

3.1 Requisiti da soddisfare in prima priorità

- Backup centralizzato dei dati di utente, con volume di dati massimo dell'ordine di 200 GB non compressi la settimana. Questo valore massimo di 200 GB/settimana è stato preso come valore di prima approssimazione in base alla situazione attuale dei servizi di backup nelle sezioni. La soluzione deve essere scalabile a volumi maggiori o minori.
- Dati distribuiti sulla LAN (anche in percentuali superiori al 90 %)
- Non vengono considerate in prima priorità le esigenze di backup di immagini di volumi o di dischi sistema, in quanto legate a soluzioni proprietarie.
- Provenienza dei dati:
 - Filesystem montati su piattaforma UNIX/Posix (Linux, Digital UNIX, Solaris, HP-UX, AIX, IRIX);
 - Filesystem esportati via NFS;
 - Filesystem visibili via AFS;
 - Filesystem montati su piattaforma Windows NT;
 - Filesystem esportati via Netbios.

- Ubicazione del server di backup: piattaforma Linux, Digital UNIX, Solaris o Windows NT.
- Destinazione dei dati:
 - area di stage su disco per il backup e il restore immediato online;
 - unità nastro SCSI singola oppure robotizzata per l'archiviazione dei dati di backup a medio termine (1-3 anni).
Non vengono considerate le esigenze di archiviazione storica a più lungo termine.
- Restore dei dati da parte del gestore del sistema:
 - con possibilità di accedere all'indice storico dei dati archiviati per individuare il file o filesystem da recuperare;
 - con possibilità di navigare nella struttura dei filesystem salvati e di ricercare per nome i file, alberi di file o filesystem da recuperare, indipendentemente dalla struttura fisica dei dati sull'unità di archiviazione e dal sistema operativo di provenienza;
 - con procedure basate su script o linea di comando;
 - con possibilità di ripristinare le ownership ed ACL originali dei singoli file.
- Backup incrementale (possibilità di trasferire a richiesta sull'unità di archiviazione i soli file modificati dopo una certa data o ora).
- Organizzazione dei log del backup che permetta un facile accesso alla storia delle modifiche di un singolo file.

3.2 Requisiti da soddisfare in seconda priorità

- Procedura di restore non assistita da parte dell'utente.
- Procedura di restore basata su interfaccia grafica.
- Controllo della dimensione dei filesystem inclusi nel sistema di backup, con soglie di allarme.
- Possibilità di creare backup immagine (blocchi fisici) di un disco, con restore del contenuto del disco su un'altra macchina.

- Possibilità di creare un set di backup da parte dell'utente in modalità non assistita ("chiosco di backup").

Capitolo 4

Valutazione di prodotti commerciali e soluzioni *public domain*

È stata quindi verificata la rispondenza ai requisiti sopra delineati (vedi Capitolo 3) per una serie di prodotti commerciali, ed in particolare:

- HP Omniback
- Legato Networker
- IBM Tivoli (ex ADSM)
- ARCServeIT

Si è cercato inoltre di valutare un nuovo prodotto (Teradactyl), le cui caratteristiche dichiarate apparivano particolarmente vicine ai requisiti sopra formulati. Non è stato tuttavia possibile ottenere una versione di prova del pacchetto, in quanto la ditta produttrice ha affermato di non poter esportare il prodotto a causa di problemi legali. Con lo stesso criterio infine è stata valutata una soluzione integrata di backup *public domain*, Amanda[14]. Nel seguito sono elencate le caratteristiche delle prove e la rispondenza ai vari requisiti.

4.1 Terminologia

Nella valutazione che segue si identificano alcune componenti del sistema integrato di backup:

Cliente di backup: Con questo termine si definiscono le macchine (host) sulle quali risiedono fisicamente i dati di backup. Tali host sono in linea di principio distribuiti su una rete locale.

Server di backup: È la macchina (host) sulla quale risiede il software incaricato di coordinare, programmare e controllare l'esecuzione delle procedure di backup che causano il trasferimento di dati dai *clienti di backup* alle unità di memorizzazione di massa (stage su disco e unità nastro). Nel caso più semplice, tali unità di memorizzazione di massa sono direttamente collegate al *server di backup*.

Device manager: Si definisce in questo modo una macchina (host), alla quale sono collegate fisicamente le unità di memorizzazione di massa (stage su disco o nastro) verso le quali vengono trasferiti i dati di backup. Alcuni dei sistemi analizzati nel seguito ammettono la configurazione di *device manager* remoti e distinti dal *server di backup*. Nel caso più semplice, il *device manager* coincide col *server di backup*.

4.2 Configurazione delle prove

4.2.1 Amanda

Versione 2.4.1p1, installata su un PC con Linux (RedHat 6.0, kernel 2.2.5-22smp) al quale è stato collegato un jukebox DLT Digital TZ875.

4.2.2 Omniback

Versione A.03.00, installata su un PC con WNT 4.0 con un jukebox Digital TZ887N.
Versione A.03.00, installata su HP-UX 10.20 con un DLT TZ87N.

4.2.3 Networker

Versione 4.4.2, installata su Digital UNIX 4.0D con un jukebox Digital TL891.

4.2.4 Tivoli

Versione 3.1 (ADSM), server su SUN Solaris 2.5.1, storage su disco, con unità DLT4000 singola. Client su SUN Solaris 2.5.1
Versione 3.1.2.40 (ADSM), server su WNT 4.0 con un jukebox Digital TZ887N.

4.2.5 ARCserveIT

Versione 6.6 enterprise edition, server WNT 4.0, con unità nastro Exabyte Mammoth 1428

4.3 Volume dei dati e scalabilità

I requisiti prevedono un volume indicativo di 200 GB la settimana scalabili.

4.3.1 Amanda

Prevede il backup diretto su nastro con un'area di staging su disco che permette la scrittura continua su nastro. È possibile gestire un database delle versioni con un software opzionale.

4.3.2 Omniback e Networker

Prevedono il backup diretto su nastro e utilizzano un'area su disco per il database delle sessioni di backup e delle versioni dei file (dimensione stimata 500MB).

4.3.3 Tivoli e ARCServeIT

La destinazione dei dati di backup, su area di stage su disco, oppure direttamente su unità nastro, è completamente configurabile dall'amministratore.

4.4 Distribuzione sulla LAN

I prodotti analizzati hanno una struttura client-server che permette una configurazione e programmazione dei backup centralizzata.

4.4.1 Amanda

Il server centrale interroga e pilota l'attività dei clienti. La sicurezza è basata su un file di controllo accesso del formato di `.rhosts`, oppure su kerberos. Il protocollo di comunicazione è privato (non vengono utilizzati servizi "r").

4.4.2 Omniback, Networker, Tivoli, ARCserveIT

Supportano tutti *clienti* di backup (host sui quali risiedono i dati di backup) distribuiti su rete, e permettono anche l'installazione di ulteriori server con unità a nastro (*device server*) integrati nel sistema in modo trasparente. Questi alcuni dettagli della comunicazione via rete:

Omniback: il server centrale lancia un demone (tramite `inetd`) sui clienti configurati nel pool di backup. Il backup può essere lanciato anche da linea di comando.

Networker, ARCServeIT: Il server centrale pilota un demone (con protocollo privato) che gira sempre sui clienti.

Tivoli: La richiesta di backup viene pilotata dal cliente, che può essere programmato dal server in anticipo.

4.5 Commenti sulla facilità di installazione di clienti e server

4.5.1 Amanda

Trattandosi di un prodotto di *public domain*, possono verificarsi problemi di installazione sulle release dei sistemi operativi più recenti. La gestione dei file di configurazione del server è piuttosto complicata, e la documentazione è piuttosto carente.

4.5.2 Omniback

L'installazione è semplice seguendo i manuali. I clienti vanno installati manualmente con uno script. Esiste un sistema di installazione dal server che non è stato provato ma che sembra funzionare solo per i clienti PC e HP-UX. C'è un tool grafico per la configurazione e la gestione. Offre il discovery dei file-system di un cliente.

4.5.3 Networker

L'installazione è semplice seguendo i manuali. I clienti vanno installati manualmente con uno script. C'è un tool grafico per la configurazione e la gestione. La configurazione può essere fatta anche via script. Anche il backup può essere lanciato da linea di comando.

4.5.4 Tivoli

L'installazione del server è semplice (con *wizard* di installazione per Windows). È stato riscontrato qualche problema di compatibilità con le librerie di sistema di Windows. I clienti vengono installati manualmente con un apposito script. Con la versione 3.7 (non disponibile al momento della prova) sono possibili installazione ed aggiornamento automatico dei clienti da parte del server di backup.

4.5.5 ARCServeIT

L'installazione del server è semplice (con *wizard* di installazione per Windows). I clienti vengono installati manualmente con un apposito script. Possibile installazione di clienti remoti Windows e UNIX dal server.

4.6 Immagini di volumi

Non sono state fatte delle prove specifiche in quanto questa è una problematica da affrontare in seconda priorità.

4.6.1 Amanda

Il sistema è in grado di gestire l'output di "dump" Unix.

4.6.2 Omniback, Networker, ARCServeIT

Prevedono l'equivalente di un dump del disco sistema. Per i PC è previsto il salvataggio della configurazione (registry).

4.6.3 Tivoli

È in grado di gestire un "dump" Unix, oppure una copia immagine di un disco di Windows, compreso il salvataggio del registry.

4.7 Provenienza dei dati

Nelle tabelle che seguono è riassunto il supporto da parte dei vari prodotti presi in esame dei filesystem indicati nei requisiti al Capitolo 3.

Filesystem	Amanda	Omniback	Networker
Filesystem locali UNIX	OK	Digital UNIX, Solaris, HP-UX, AIX, IRIX, altri (SCO, etc.) Linux entro 1999	Linux, Digital UNIX, Solaris (anche Intel), HP-UX, AIX, IRIX, altri (SCO, UnixWare, etc)
Filesystem visibili via NFS	OK	OK	OK
Filesystem visibili via AFS	solo file, non ACL	solo file, non ACL	solo file, non ACL
Filesystem locali WNT	NO	WNT W98 W95	WNT W98 W95
Filesystem visibili via Netbios	Dati leggibili via SAMBA. No ACL.	?	?
Filesystem locali MAC	NO	NO	OK
Database vari (Oracle, Informix, SAP, Microsoft SQL, ecc.)	NO	OK (lic. aggiuntiva)	OK (lic. aggiuntiva)

Filesystem	Tivoli	ARCServeIT
Filesystem locali UNIX	Digital UNIX, Solaris, HP-UX, AIX, IRIX, Linux (senza supporto IBM), altri (SCO, etc.)	Digital UNIX, Solaris HP-UX, AIX, IRIX, Linux (beta) altri (SCO, etc.)
Filesystem visibili via NFS	OK	NO (problema di licenze)
Filesystem visibili via AFS	OK (solo su AIX)	OK (licenza aggiuntiva)
Filesystem locali WNT	WNT W98 W95	WNT W98 W95
Filesystem visibili via Netbios	OK	NO
Filesystem locali MAC	OK	OK
Database vari (Oracle, Informix, SAP, Microsoft SQL, ecc.)	OK (lic. aggiuntiva)	OK (lic. aggiuntiva)

4.8 Ubicazione del server

4.8.1 Amanda

Linux, Digital UNIX, Solaris, e dovunque possano essere compilati i prodotti GNU. La distribuzione utilizzata nella prova su Linux ha richiesto una modifica a causa di una variazione nella struttura degli include file in RedHat 6.0.

4.8.2 Omniback

HP-UX, WNT, Solaris (solo *device manager*), AIX (solo *device manager*), IRIX (solo *device manager*). Il *device manager* è l'host al quale sono collegati fisicamente i dispositivi di backup. È necessario almeno un server completo (*cell manager*) per ogni sito.

4.8.3 Networker

Digital UNIX, Solaris, HP-UX, AIX, IRIX, WNT.

4.8.4 Tivoli

AIX, Solaris, HP-UX, WNT, OS400.

4.8.5 ARCServeIT

AIX, Solaris, HP-UX, WNT, Linux (beta).

4.9 Destinazione dei dati

Lo stage su disco viene gestito privatamente dal server. Le unità a nastro singolo SCSI sono gestite con i driver standard del Sistema Operativo. Ci sono delle differenze per i jukebox di nastri.

4.9.1 Amanda

Viene supportato chio su HP-UX, ma non su Solaris, o Digital UNIX. Per la prova su Linux ha dovuto essere aggiunto un programma di controllo del jukebox (una versione molto semplificata del pacchetto “juke” di Fermilab). È necessario un po’ di sviluppo per un supporto estensivo di questa funzione.

4.9.2 Omniback

Gestisce i jukebox SCSI2 dove la meccanica ha lo stesso ID del driver ma LUN diverso. Nella versione per HP-UX bisogna creare un device apposito per la meccanica del jukebox. Questa operazione alle volte è complicata. Nella versione WNT c’è un tool che riconosce e configura automaticamente il jukebox. Supporta anche jukebox di magneto-ottici. Si possono creare pool di device: quando viene richiesto un media lo si può montare in un qualsiasi device del pool. È possibile definire gli slot da utilizzare e quelli da lasciare liberi per altre applicazioni.

4.9.3 Networker, Tivoli, ARCServeIT

Gestisce i jukebox SCSI2 dove la meccanica ha lo stesso ID del driver ma LUN diverso. Bisogna creare un device apposito per la meccanica. Supporta anche jukebox di magneto-ottici. Si possono creare pool di device: quando viene richiesto un media lo si può montare in un qualsiasi device del pool. È possibile definire gli slot da utilizzare e quelli da lasciare liberi per altre applicazioni.

4.10 Backup dei dati

4.10.1 Amanda

Il backup incrementale è supportato, anche se organizzato in modo dinamico a seconda delle risorse disponibili. Esiste infatti uno schema di priorità che può eventualmente far sì che il backup di un certo set di dati venga rimandato se non vi sono risorse di rete o di storage disponibili.

4.10.2 Omniback, Networker, Tivoli, ARCServeIT

I dati possono essere compressi dai clienti riducendo il traffico sulla rete. Inoltre si possono definire comandi particolari da eseguire prima del backup e dopo il restore. È possibile combinare un particolare tipo di scheduling (es. mensile, settimanale, giorno per giorno) col tipo di backup desiderato (full, incrementale a livelli). Ci sono varie tipologie di segnalazioni di eventi particolari (es. report del backup, media-full, client not available, ecc.). ARCServeIT e Networker comprendono anche la segnalazione via e-mail. Omniback e Networker dispongono di una funzione di preview del backup che simula il backup, contatta i clienti coinvolti e segnala possibili malfunzionamenti; Omniback e Tivoli indicano le dimensioni dei dati da salvare e stimano il rapporto di compressione. Il backup si può lanciare anche via script.

4.11 Restore dei dati e accesso alla storia dei file

4.11.1 Amanda

Indice storico: È incluso un pacchetto (indipendente dal nucleo del pacchetto di backup) che gestisce un indice di tutti i dati archiviati con la data di archiviazione e lo rende disponibile attraverso un apposito server di rete.

Navigazione nella struttura dei filesystem salvati: OK, utilizzando l'utility "amrecover".

Script e linea di comando: negli script e nella linea di comando può essere utilizzato il comando esplicito di restore di file e directory (amrestore).

Ripristino di ownership ed ACL: senza problemi per UNIX. Il restore dei file avviene su UNIX per i file esportati via SAMBA, per cui è necessaria una traduzione degli attributi. Nessun supporto per AFS.

Accesso alla storia dei file via log: questa funzione è limitatamente possibile nell'utility "amrecover". Il server degli indici salva tuttavia solo la data di backup ed il nome dei file salvati, e non la loro ultima data di modifica o la loro dimensione.

4.11.2 Omniback, Networker, Tivoli e ARCServeIT

Con l'interfaccia grafica è possibile navigare attraverso la struttura del file-system salvato e selezionare i file o i sotto-alberi da recuperare. È possibile selezionare la versione del file che interessa. In Tivoli la modalità di conservazione delle versioni dipende dalla *policy* di backup selezionata. L'indice delle versioni di Omniback è per variazione del file (ad esempio se il file è cambiato 3 volte in 20 giorni ci sono 3 versioni nell'indice) e non

per data di backup come su Networker (con lo stesso esempio di Omniback ci sarebbero 20 entry nell'indice) e quindi risulta più semplice e immediato. Viene garantita la riservatezza: si possono definire utenti e autorizzarli o meno a compiere specifiche azioni. Si può lanciare il restore anche tramite una serie di comandi contenuti in uno script. Tivoli possiede un log dettagliato, file per file. ARCServeIT possiede un log dettagliato e configurabile. Omniback e Networker producono report dettagliati (particolarmente dettagliati quelli di Omniback) dell'attività, ma non in formato di testo leggibile. Tutti i prodotti comprendono procedure di *disaster recovery*.

4.12 Formato dei dati su nastro e formato database

4.12.1 Amanda

Tutti i nastri utilizzati dal sistema devono essere preventivamente etichettati. I nastri contengono, dopo l'etichetta, una serie di saveset separati da filemark. I saveset contengono dati raccolti con dump o gnu-tar, con un record iniziale aggiunto, e scritti in record fisici di 32k. I file possono essere opzionalmente compressi (gzip) dal cliente o dal server prima di essere scritti su nastro. Volendoli rileggere senza utilizzare amanda, dopo essersi posizionati sul file corretto, occorre utilizzare un comando del tipo:

```
dd if=/dev/tape bs=32k skip=1 [| gzip -d] | gtar -xvf -.
```

4.12.2 Omniback, Networker, Tivoli e ARCServeIT

I nastri devono essere inizializzati con una label che li identifica univocamente nel database. È previsto un "template" per la numerazione automatica della label dei media (es. gruppo_NNN dove NNN è un numero progressivo). Esiste un database che contiene tutte le informazioni sui backup eseguiti e l'indice dei file salvati divisi per pool di backup. Si specifica dopo quanto tempo il media può essere sovrascritto e per quanto tempo si deve mantenere l'indice dei file in esso contenuti. C'è la possibilità di identificare le cassette con un codice a barre. Il formato del database e dei nastri è proprietario. Esistono comandi per la manutenzione del database.

4.13 Condizioni di commercializzazione

Analizziamo qui le offerte dei prodotti commerciali pervenute al gruppo di lavoro. Le offerte sono formulate per un'applicazione sulla scala delle 31 sedi INFN. Tutti i costi indicati si intendono i.v.a. esclusa.

4.13.1 Omniback

Server Unix (SOLO HP-UX). Comprendente un numero illimitato di clienti eterogenei, autoloader con < 200 slot e supporto per *un* drive nastro: 6.23ML. Per realizzare la funzione di *device manager* è necessario acquistare un'altra licenza server.

Server WNT alle stesse condizioni del server Unix: 1.75ML.

Licenza per poter utilizzare più di un drive sullo stesso server: (multidrive) 3.36 ML.

Nota: questi prezzi, a differenza degli altri, sono stati offerti ad una singola sezione, e sono suscettibili di un trattamento di maggior favore per l'acquisto centralizzato INFN.

4.13.2 Networker

Server Unix: (con 10 Clienti UNIX) 12.1ML + 3 ML per poter utilizzare clienti Windows (compresa maintenance per un anno).

Server WNT: (con 10 clienti WNT) 3.3 ML + 3 ML per poter utilizzare clienti Unix

Licenze per clienti aggiuntivi: 2.5 ML per 5 clienti, 10 ML per 25 clienti, 30 ML per 100 clienti.

Licenze per poter utilizzare autoloader: 3 ML se < 8 slot, 7 ML se < 16 slot, 10 ML se < 32 slot, 13 ML se < 64 slot.

4.13.3 Tivoli

Server Unix (AIX o Solaris): 12 ML con "abilitazione di rete" (solo server).

Server WNT: 5.8 ML con "abilitazione di rete" (solo server).

Licenze clienti: cliente singolo: 130 KL, 50 clienti 6.2 ML, 100 clienti 11.6 ML

Licenze per librerie di nastri non coperte dalla licenza base: 17.4 ML per server Unix, 14.9 ML per server WNT.

4.13.4 ARCServeIT

Il prodotto viene distribuito in varie versioni. È stata considerata la versione "Advanced edition", il cui prezzo di listino è di 4.1 ML per server WNT con libreria di 7 nastri e 20 clienti WNT, più 621 KL per ogni cliente Unix. Si considera che questa soluzione sia svantaggiosa per la maggior parte delle sezioni, che hanno un gran numero di clienti di

backup installati. È stata quindi considerata la versione “enterprise” del prodotto, per la quale il distributore italiano, Computer Associates, si è dimostrato riluttante nel produrre un’offerta per il solo prodotto di backup, in ragione del fatto che esiste un complicato meccanismo di licensing basato su una stima del “peso”, in termini sostanzialmente di potenza di CPU e di capacità disco, degli host (clienti di backup) coinvolti. È stata richiesta dunque un’offerta per una configurazione “tipo” corrispondente a 24 macchine Unix e 22 PC con Windows (con caratteristiche hardware corrispondenti alle macchine della sezione INFN di Padova). L’offerta commerciale risultante ammonta a 82.7 K\$, corrispondenti a 150 ML circa.

In base ai risultati ottenuti, il gruppo di lavoro ha segnalato alla Commissione Calcolo dell’INFN i prodotti nel seguente ordine di preferenza: Tivoli, Omniback, Legato, ArcServeIT. La Commissione, nella seduta del 4 Marzo 2000 ha optato per la scelta che offre la maggiore convenienza economica, e cioè Omniback.

	Omniback	Networker	Tivoli	ARCserveIT
Destinazione dei backup	Diretto su nastro + area su disco per database delle sessioni di backup e delle versioni dei file	Diretto su nastro + area su disco per database delle sessioni di backup e delle versioni dei file	Configurabile: stage su disco o diretto su nastro	Configurabile: stage su disco o diretto su nastro
Distribuzione su LAN	Il server lancia un demone via inetd sui clienti	Il server pilota un demone che gira sui clienti	La richiesta parte dal cliente, programmabile in anticipo dal server	Il server pilota un demone che gira sui clienti
Installazione	Semplice	Semplice	Semplice (è stato riscontrato qualche problema con le librerie di Windows)	Semplice
Ubicazione server	HPUX, WNT, Solaris ¹ , AIX ¹ , IRIX ¹	Digital Unix, HPUX, Solaris, AIX, IRIX, WNT	AIX, Solaris, HPUX, WNT, AS400	AIX, Solaris, HPUX, WNT, Linux (beta)
Provenienza dei dati				
Immagini volumi Unix	OK	OK	OK	OK
Dump registry di NT	OK	OK	OK	OK
Supporto filesystem locali Unix	Digital UNIX, Solaris, HP-UX, AIX, IRIX, altri (SCO, etc.) Linux entro 1999	Linux, Digital UNIX, Solaris (anche Intel), HP-UX, AIX, IRIX, altri (SCO, UnixWare, etc)	Digital UNIX, Solaris, HP-UX, AIX, IRIX, Linux (senza supporto IBM), altri (SCO, etc.)	Digital UNIX, Solaris, HP-UX, AIX, IRIX, Linux (beta) altri (SCO, etc.)
Supporto filesystem visibili via NFS	OK	OK	OK	NO (problema di licenze)
Supporto filesystem visibili via AFS	solo file, non ACL	solo file, non ACL	OK (solo su AIX)	OK (licenza aggiuntiva)
Supporto filesystem locali WNT	WNT W98 W95	WNT W98 W95	WNT W98 W95	WNT W98 W95
Supporto filesystem visibili via Netbios	?	?	OK	NO
Supporto filesystem locali Mac	NO	OK	OK	OK
Supporto database vari (Oracle, Informix, SAP, Microsoft SQL, ecc.)	OK (lic. aggiuntiva)	OK (lic. aggiuntiva)	OK (lic. aggiuntiva)	OK (lic. aggiuntiva)
Destinazione dei dati				
Gestione jukebox di nastri SCSI2	OK	OK	OK	OK
Gestione jukebox di magneto-ottici	OK	OK	OK	OK
Gestione pool di device	OK	OK	OK	OK

¹Solo device manager

Tabella 4.1: Tabella riassuntiva delle funzionalità dei prodotti commerciali esaminati (prima parte).

	Omniback	Networker	Tivoli	ARCserveIT
Backup dei dati				
Compressione client-side	OK	OK	OK	OK
Supporto comandi pre-backup e post-restore	OK	OK	OK	OK
Segnalazione eventi particolari	OK	OK <small>anche via e-mail</small>	OK	OK <small>anche via e-mail</small>
Preview del backup	OK	OK	NO	NO
Stima compressione	OK	NO	OK	NO
Restore dei dati e accesso alla storia dei file				
Navigazione (grafica) nel filesystem salvato	OK	OK	OK	OK
Restore di una certa versione di un file	OK <small>indice per variazione</small>	OK <small>indice per data</small>	OK	OK
Esecuzione del restore tramite script	OK	OK	OK	OK
Log files	OK (non in formato testo)	OK (non in formato testo)	OK	OK
Disaster recovery	OK	OK	OK	OK
Formato dei dati su nastro e formato database				
Inizializzazione dei nastri con label	OK	OK	OK	OK
Identificazione nastri con codice a barre	OK	OK	OK	OK
Formato dei nastri	Proprietario	Proprietario	Proprietario	Proprietario
Formato database	Proprietario	Proprietario	Proprietario	Proprietario

Tabella 4.2: Tabella riassuntiva delle funzionalità dei prodotti commerciali esaminati (seconda parte).

	Omniback	Networker	Tivoli	ARCserveIT
Condizioni di commercializzazione				
Server Unix	6.23ML ^{2 3} solo HPUX	(compresi 10 clienti Unix) 12.1ML + 3ML per poter utilizzare clienti Windows	12 ML solo Solaris o AIX con "abilitaz. di rete"	
Server WNT	1.75 ML ^{2 3}	(compresi 10 clienti WNT) 3.3ML + 3ML per poter utilizzare clienti Unix	5.8ML con "abilitaz. di rete"	4.1ML⁴
Clienti	Compresi nella licenza del server	5 cl.: 2.5ML 25 cl.: 10ML 100 cl.: 30ML	1 cl.: 130 KL 50 cl.: 6.2ML 100 cl.: 11.6ML	621KL⁴ per ogni client Unix
Altre licenze	Multidrive sullo stesso server: 3.36ML ³	Autoloader: 3ML se <8 slot, 7ML se <16 slot, 10ML se <32 slot, 13ML se <64 slot	Librerie di nastri non comprese nella lic. base: 17.4ML per Unix, 14.9 per WNT	

²Comprende un numero illimitato di clienti eterogenei, autoloader con max 200 slot e un drive nastro. Per realizzare la funzione di device manager è necessario acquistare un'altra licenza server.

³Questi prezzi, a differenza degli altri, sono stati offerti ad una singola sezione, e sono suscettibili di un trattamento di maggior favore per l'acquisto centralizzato INFN.

⁴Advanced Edition

Tabella 4.3: Riassunto delle condizioni di commercializzazione dei prodotti commerciali esaminati.

Capitolo 5

Technology tracking di soluzioni e supporti magnetici.

Vengono prese in esame le tecnologie compatibili con il volume di dati identificato nei requisiti in precedenza descritti (vedi Capitolo 3). Non verranno quindi presi in esame tutti quei sistemi di archiviazione come ad esempio i CD (*Compact Disk*) poiché al momento di questa analisi i costi e le capacità di tali sistemi non risultano essere competitivi con quelle di tipo nastro.

5.1 Le tecnologie nastro

L'unica tecnologia di archiviazione disponibile con costi e caratteristiche accettabili è quella a nastro. Nell'ambito di questa tecnologia si distinguono due tipologie, lineare ed elicoidale, che si differenziano per il modo di interfacciamento meccanico tra la testina ed il nastro e per la struttura con cui vengono realizzati i nastri magnetici[9].

Nel caso dei nastri lineari è il solo nastro magnetico a scorrere rispetto alla testina magnetica e le tracce sono parallele alla direzione di trascinamento del nastro stesso. Nel caso dei sistemi a nastro di tipo elicoidale sia il nastro che la testina si muovono reciprocamente con la seconda che ruota su un asse obliquo rispetto alla direzione di trascinamento del nastro. Questo consente a parità di testina magnetica di incrementare la densità di scrittura, vedi eq.5.1.

$$Den_{stape} = \frac{Den_{thead}}{\cos(\Theta_{thead})} \quad (5.1)$$

Nelle unità nastro di tipo lineare la superficie di contatto tra testina e nastro è molto ridotta, al contrario di quelle elicoidali.

Questi due differenti approcci determinano differenti caratteristiche dei diversi sottosistemi. In particolare si avrà una maggiore densità di scrittura per i sistemi di tipo

elicoidale ma una minore durata sia per quello che riguarda il nastro che per la testina di scrittura/lettura.

Chiaramente l'affidabilità dei sistemi a nastro lineare è ben diversa da quella dei sistemi elicoidali, sia in termini di deterioramento dell'unità nastro che del nastro magnetico.

Il divario di prestazioni, sia per quanto concerne la velocità di lettura/scrittura, che per la capacità totale, si è recentemente ridotto, ferma restando la maggiore affidabilità dei prodotti a nastro lineari.

Ovviamente anche la tecnologia lineare ha delle limitazioni dovute agli attriti di tipo meccanico. Come vedremo, alcuni costruttori, al fine di accrescere l'affidabilità dei loro prodotti, hanno adottato soluzioni che però hanno fatto crescere notevolmente i costi del singolo drive, il quale peraltro raggiunge dimensioni considerevoli.

Modello	Capacità Nativa in (GB)	Formato di Scrittura	Compatibilità
DLT	10, 15, 20, 35, 40	Lineare	Totale
SLR	12, 16, 25	Lineare	Parziale
STK 9840	20	Lineare	Assente
Magstar	10, 20	Lineare	Parziale
MagstarMP	5	Lineare	Assente
DAT	1.3, 2, 4, 12, 20	Elicoidale	Totale
AIT	25, 50	Elicoidale	Totale
Mammoth	20, 40	Elicoidale	Parziale

Tabella 5.1: Lista dei modelli di nastro presi in esame

In tabella 5.1 sono riepilogate le differenti unità nastro prese in esame. Alcune linee di prodotto sono da lungo tempo sul mercato con caratteristiche sempre migliori, come DAT e DLT. Tali linee di prodotto hanno mantenuto, nel tempo, la compatibilità con i prodotti precedenti.

Altri vendor hanno sviluppato soluzioni non completamente compatibili con le loro precedenti linee di prodotto: Magstar, Mammoth, SLR. In qualche caso (vedi MagstarMP) all'interno della medesima tecnologia sono state realizzate due linee distinte di prodotto meccanicamente incompatibili. Inoltre, MagstarMP è disponibile, data la sua capacità per singola cassetta, solo nella soluzione libreria.

5.2 Criteri di confronto delle tecnologie nastro

Al fine di poter confrontare le diverse tecnologie è necessario individuare una serie di parametri caratteristici oggettivi tali da consentire una corretta valutazione delle differenti soluzioni.

Sono state individuate le seguenti caratteristiche principali:

Capacità S'intende la capacità del nastro in formato nativo, senza l'applicazione di nessun meccanismo di compressione, poichè quest'ultima dipende, oltre che dall'algoritmo implementato dal costruttore, anche dalla tipologia dei dati da immagazzinare.

Velocità S'intende la velocità con cui si può accedere in lettura/scrittura al nastro in modo sequenziale e sostenuto, sempre riferendosi alla scrittura in formato nativo.

Ricerca file: tempo medio di ricerca di un file su nastro.

Affidabilità Solitamente il parametro che caratterizza questo aspetto è l'**MTBF** (*Mean Time Between Failures*), corrispondente al tempo medio tra due errori dell'elettronica presente nel drive. In realtà questa informazione da sola può indurre in errore, infatti prodotti molto diversi hanno un **MTBF** simile tra di loro.

A questo proposito risulta indispensabile conoscere un altro parametro, che spesso non viene indicato, noto come *Duty Cycle*, il quale corrisponde al tempo medio di utilizzo, che nel caso dei nastri elicoidali è sempre inferiore al 30%.

Diventa perciò indispensabile pensare a qualche cosa di differente. Un lavoro dell'HP [2] al fine di poter confrontare i differenti prodotti, definisce un nuovo parametro "*Predicted MTBF*" nel modo seguente:

$$\text{Predicted MTBF} = \frac{\text{MTBF} * \text{DutyCycle}}{\text{UserDutyCycle}} \quad (5.2)$$

dove *User Duty Cycle* non è altro che la percentuale di tempo in cui effettivamente viene impegnato il drive¹.

Un altro parametro essenziale che determina l'affidabilità è dato dalla durata della testina magnetica (*Head Life*) e dalla durata del supporto magnetico. Quest'ultima va intesa sia come durata nel tempo (*archivio storico*) che per le successive riscritture (*politiche di riciclaggio dei nastri*).

¹La percentuale nuova non dovrà essere molto diversa da quella originale in quanto si rischia di lavorare fuori dalle specifiche

Quest'ultimo punto va tenuto in considerazione al fine di ridurre i costi di esercizio del sistema di Backup, attuando una opportuna politica di riutilizzo delle cassette.

Compatibilità Indica la possibilità di utilizzare il drive di ultima generazione al fine di poter accedere sia in lettura e possibilmente anche in scrittura le cassette prodotte, in altre sedi, con modelli precedenti.

Questo fatto deve essere tenuto in considerazione in quanto consente di proteggere l'investimento e non costringe a mantenere tecnologie differenti allo stesso tempo.

Write Format Indica il formato con il quale viene scritto fisicamente il nastro: lineare o elicoidale.

Integrazione Riguarda la possibilità di disporre delle unità nastro all'interno di differenti sistemi di libreria.

Costo Deve prendere in considerazione sia il costo dei singoli drive (**Costo Drive**) che quello di archiviazione (**Costo x MB**).

5.3 Confronti

Riepilogo dei costi dei drive e delle cassette i.v.a. esclusa :

Modello Drive	Costo Drive (MI)	Cassetta	Cassetta Pulizia	Costo x MB
DLT 35GB	10.	150,000	110,000	4.2
SLR 12GB	2.8	289,000	88,000	23.5
DAT 12GB	2.0	45,000	13,000	3.7
AIT 25GB		145,000		5.7
Mammoth 20GB	7.5	150,000	42,000	7.3
Magstar 20GB	30.	150,000	110,000	7.3
STK 20GB	30.	150,000	110,000	7.3

Tabella 5.2: Costi dei drive e dei media

Questa è la lista dei principali prodotti attualmente disponibili sul mercato:

5.4 Librerie Nastro

Al contrario di quello che accade per le unità nastro, essenzialmente patrimonio di poche società, vi è un elevatissimo numero di società che di fatto integra all'interno delle proprie librerie i più diffusi sistemi a nastro.

Tabella 5.3: Caratteristiche dell'unità nastro DLT della Quantum[1]

<i>Duty Cycle: 100% / FMT: Lineare / Compatibilità: Totale</i>					
Modello	Cap.	Speed(MB/s)	MTBF(h)	HL(h)	Ric.File.(s)
8000	40	6.0	400,000	30,000	60.
7000	35	5.0	200,000	30,000	60.
4000	20	1.5	200,000	10,000	68.
2000XT	15	1.25	80,000	30,000	68.

Tabella 5.4: Caratteristiche dell'unità nastro AIT di Sony [5]

<i>Duty Cycle: 40% / FMT: Elicoidale / Compatibilità: Totale</i>					
Modello	Cap.	Speed(MB/s)	MTBF(h)	HL(h)	Ric. File.(s)
2	50	6.0	200,000	30,000	37.
1	35	3.0	200,000	30,000	37.

Tabella 5.5: Caratteristiche dell'unità nastro SLR di Tandberg[4]

<i>Duty Cycle: 10% / FMT: Lineare / Compatibilità: Limitata</i>					
Modello	Cap.	Speed(MB/s)	MTBF(h)	HL(h)	Ric. File.(s)
50	25	2.0 - 1.0	300,000	??????	30.
32	16	1.5 - 0.75	300,000	??????	30.
24	12	1.2 - 0.6	300,000	??????	30.

Tabella 5.6: Caratteristiche dell'unità nastro DAT di HP[3]

<i>Duty Cycle: 12% / FMT: Elicoidale / Compatibilità: Totale</i>					
Modello	Cap.	Speed(MB/s)	MTBF(h)	HL(h)	Ric. File.(s)
DDS4	20	<i>Annunciato</i>			
DDS3	12	1.0	300,000	??????	50.
DDS2	4	0.512	300,000	??????	50.
DDS1	2	0.180	150,000	??????	50.

Tabella 5.7: Caratteristiche dell'unità nastro Mammoth di Exabyte [6]

<i>Duty Cycle: 12% / FMT: Elicoidale / Compatibilità: Limitata</i>					
Modello	Cap.	Speed(MB/s)	MTBF(h)	HL(h)	Ric. File.(s)
2	40	6	250,000	35,000	??
1	20	3	250,000	35,000	??
LT	14	2	250,000	35,000	??

Tabella 5.8: Caratteristiche dell'unità nastro Magstar di IBM[7]

<i>Duty Cycle: 12% / FMT: Lineare / Compatibilità: Parziale</i>					
Modello	Cap.	Speed(MB/s)	MTBF(h)	HL(h)	Ric. File.(s)
E	20	14	400,000	??????	50.
B	10	9	300,000	??????	50.
MP	5	7	300,000	??????	50.

Tabella 5.9: Caratteristiche dell'unità nastro STK 9840 di StorageTek[8]

<i>Duty Cycle: 12% / FMT: Lineare / Compatibilità: Assente</i>					
Modello	Cap.	Speed(MB/s)	MTBF(h)	HL(h)	Ric. File.(s)
9840	20	10	400,000	??????	10

Le diverse librerie si distinguono oltre che per dimensione ed espandibilità, anche per funzionalità. In linea generale si può fare una sommaria distinzione tra due categorie. In particolare vi sono, limitatamente a quelle di piccola taglia (≤ 10 cassette), gli autoloader e le librerie vere e proprie. La differenza principale è dovuta al fatto che i primi non posseggono un meccanismo di riconoscimento della cassetta (mediante lettura ottica) ed in alcuni casi non consentono di accedere in modalità randomica i vari nastri.

Tali sistemi sono di fatto molto limitati come dimensione e non sono espandibili, con conseguente impatto sui costi.

5.5 Evoluzione Futura

La situazione vede l'evoluzione di quelle che sono le odierne tecnologie; ad esempio il DLT8000 e il DDS4 hanno ottimizzato aspetti differenti che hanno accresciuto la capacità delle corrispondenti classi di sistemi a nastro; Magstar a breve porterà la capacità degli attuali modelli Magstar-E dagli odierni 20GB a 40GB semplicemente raddoppiando la lunghezza del nastro, mantenendo invariate tutte le altre caratteristiche.

Sul lungo periodo, che possiamo individuare tra 1-3 anni da oggi, verranno messe in produzione alcune tecnologie lineari attualmente in fase di definizione.

La tecnologia che è attualmente alla base del DLT ha ormai 5 anni e non ha margini di crescita ulteriori, tali perlomeno da raggiungere le centinaia di GB, le quali potranno essere raggiunte solo cambiando sia il nastro magnetico che l'ottica della testina[10]. Le principali aziende del settore come Seagate[12], HP[3] e IBM[13], al fine di accrescere le caratteristiche generali dei nastri stanno introducendo nel mercato una nuova tecnologia nota come LTO (*Linear Tape Open*). Tale tecnologia fonda le sue basi su alcune delle principali caratteristiche presenti nella linea di prodotto DLT, andando a migliorare quelli che sono gli attuali limiti tecnologici di tale prodotto. In questo modo gli sviluppi dei sistemi basati su LTO dovranno coprire come capacità l'intervallo compreso tra 100 GB, nella prima versione e 500GB nella massima evoluzione, con velocità di accesso comprese tra i 10 MB/s ed i 40 MB/s.

Una ulteriore diversificazione del prodotto verrà realizzata sulla base di due distinte richieste capacità totale e velocità di accesso. Tali linee di prodotto sono oggi denominate *Ultrium* e *Accelis*. La prima sarà caratterizzata dalla elevata capacità e dalla affidabilità il tutto utilizzando un formato del nastro confrontabile con l'odierno DLT. La seconda avrà un formato simile all'8mm ed avrà una capacità ridotta (circa 1/4).

5.6 Conclusioni

Il risultato più evidente di questo lavoro è dato dalla forte evoluzione che vi è stata da parte delle tecnologie di registrazione lineare su nastro che hanno di fatto raggiunto quelle di tipo elicoidale.

Inoltre è evidente come tale linea di sviluppo sarà mantenuta anche nei prossimi anni grazie allo sviluppo fortemente voluto dalle principali industrie di Information Technology presenti sul mercato. Le ragioni di ciò sono da ricondurre a due caratteristiche:

- Elevata affidabilità
- Costi di esercizio contenuti

Tali caratteristiche individuano nella tecnologia lineare la sola perseguibile per gli scopi che il gruppo di lavoro INFN-TOOLS si è dato (backup). Tra queste è certamente il DLT a fornire il migliore compromesso tra affidabilità, capacità e costi.

È risultato evidente come la sua linea di sviluppo nota come *SuperDLT* sia di sicuro interesse sia per le prestazioni generali che per la compatibilità, analogamente a quanto è accaduto in passato.

Quest'ultimo fatto diventa importante al fine di ridurre il numero totale di drive da supportare nel tempo, a seguito delle evoluzioni tecnologiche.

Un ulteriore argomento a favore di questa scelta è connesso alla loro integrazione nei sistemi a libreria. In questo lavoro si è ritenuta interessante la soluzione libreria Compaq per la media e grande sede INFN, in quanto oltre a fornire una buona capacità (350 GB nativi) consente di espandere, con l'adozione di opportuni box, il sistema fino a 46 cassette (1.6 TB nativi) ed un totale di 6 drive.

Il costo della soluzione base con 10 slot ed un drive 7000 è di circa 20MI, mentre l'unità di espansione di 16 cassette ne costa circa 8.

In tal senso le soluzioni alternative basate per altro su tecnologie elicoidali, come FreeFrog basato su AIT-1, presentano lo svantaggio di non fornire alcuna espandibilità e dunque non consentono la protezione dell'investimento iniziale.

Per quello che riguarda le piccole sedi, si può adottare la tecnologia DAT nel formato DDS3 o DDS4 unitamente ad un mini-autoloader con 6-8 cassette che consente una capacità totale nativa di 72 o 120 GB. Il costo di questa soluzione è di circa 4.4 MI.

5.7 Tecnologia disco

L'attuale tecnologia nastro è stata di fatto sollecitata dalla crescente capacità dei dischi, fortemente richiesta dalle più recenti applicazioni multimediali come il "video on de-

mand”. A supporto di queste accresciute esigenze sono stati sviluppati i più recenti protocolli di accesso ai dischi come SCSI, SCSI-2, SCSI-3, Ultra-SCSI, Ultra2-SCSI e più recentemente Ultra3-SCSI. Questa spinta tecnologica ha portato benefici anche sul “mercato consumer” che oggi ha la possibilità di utilizzare a basso costo (2-3 volte meno di SCSI) dischi ad elevata capacità e con prestazioni di tutto rilievo.

Dischi a 7200 giri, con interfacciamento EIDE UltraATA/66 possono essere scritti/letti a 12-13 MB/s. Le sole limitazioni sono:

- limitata espandibilità (Max 4 device)
- throughput max (12MB/s)
- latenza doppia
- prestazioni limitate in contesa di risorse

In Fig.5.7 viene mostrato il costo per GB, dei sottosistemi disco, in funzione delle differenti tecnologie di interfacciamento con il sistema. Si può notare come la tecnologia UltraDMA 66 abbia un costo 2.5 volte inferiore a quella SCSI. Per quello che riguarda le limitazioni dovute al numero complessivo di dischi che è possibile configurare, si può sempre adottare dei controller EIDE anche in versione Raid² per connettere 4 dischi per ogni controller. Ogni controller può connettere fino a 140GB ad un costo complessivo di 2,6 MLire controller compreso.

²Il Raid è disponibile solo nelle modalità 0,1

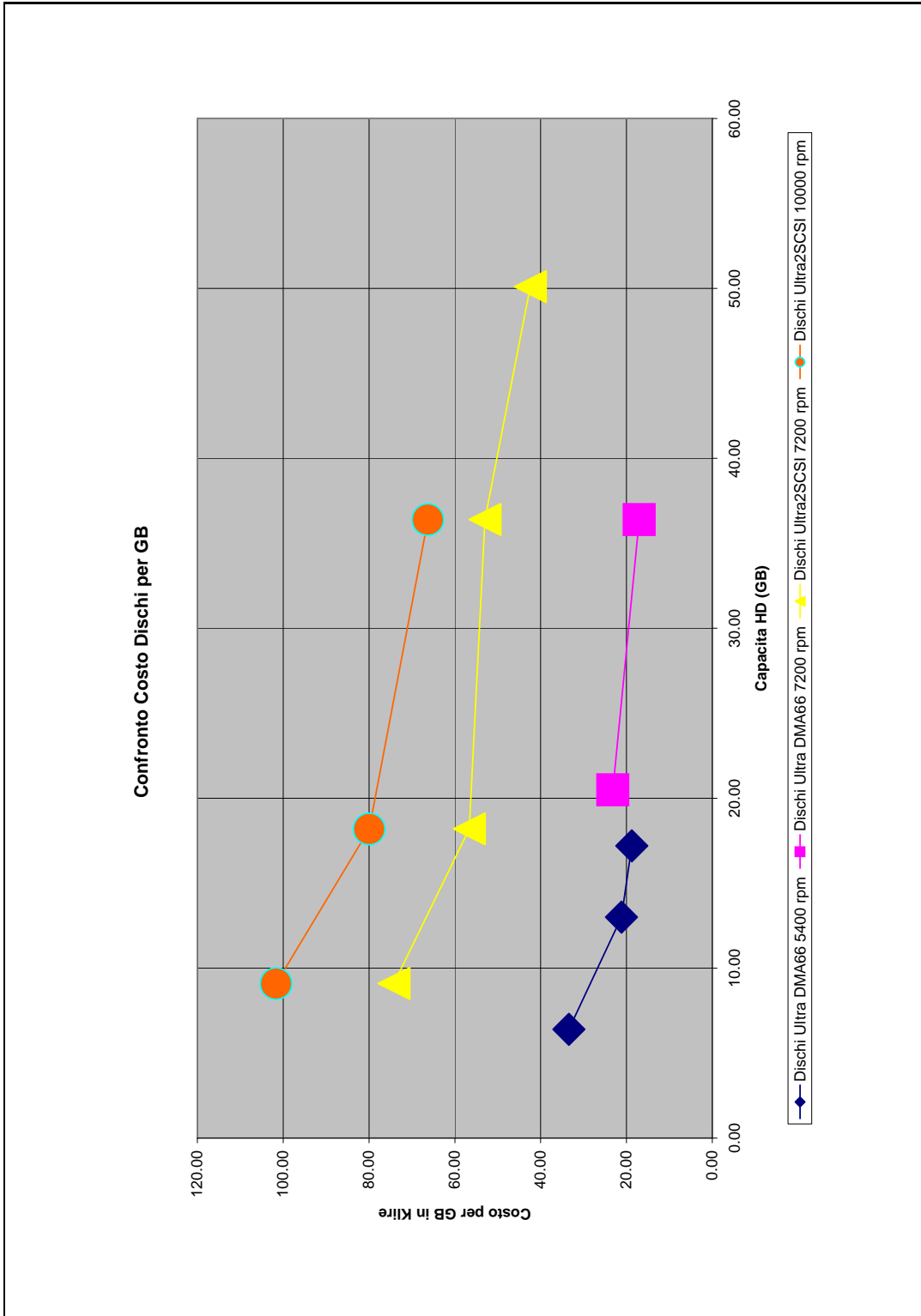


Figura 5.1: Costo dei sottosistemi disco per GB

Bibliografia

- [1] <http://www.quantum.com/>
- [2] <http://www.hp.com/tape/papers/mtbf.html>
- [3] <http://www.hp.com/tape/index.html>
- [4] <http://www.tandberg.com/>
- [5] <http://www.sony.com/>
- [6] <http://www.exabyte.com/>
- [7] <http://www.storage.ibm.com/>
- [8] <http://www.storagetek.com/>
- [9] <http://www.wwpi.com/>
- [10] <http://www.westworldproductions.com/archive/1999/0599ctr/6059.htm>
- [11] <http://www.lto.org/>
- [12] <http://www.seagate.com/support/tape/lto/ltopages.html>
- [13] <http://www.ibm.com/hardware/tape/lto/lto.html>
- [14] <http://www.amanda.org/>