

Deduplizierung mit Bacula Base Jobs

Bacula Base Jobs

Who? Philipp Storz

From? <http://www.dass-it.de/>

When? FrOSCon.de, 26.08.2012

Rev : 17207

Philipp Storz

- seit 1998 Beschäftigung mit Linux
- seit 2001 Vollzeit
- seit 2004 - Geschäftsführer dass IT GmbH
- seit 2008 - Jährliche Veranstaltung Bacula Konferenz
- 2012 Bacula Buch - Open Source Press

Eckdaten

- Geschäftsfelder rund um Open Source:
 - Consulting
 - Support
 - Maßgeschneiderte Anpassungen und Entwicklung
- Gegründet: 2004
- Mitarbeiter: 8



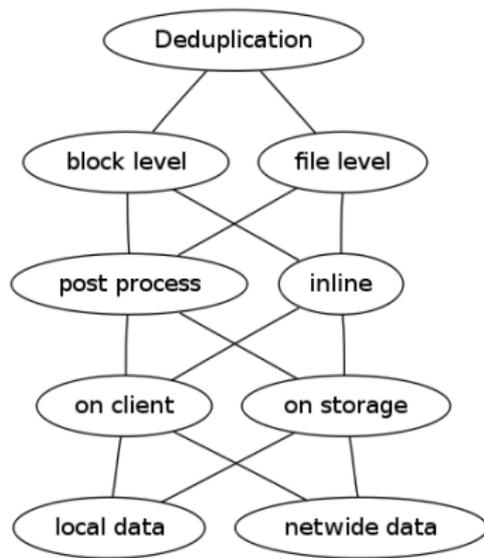
dass IT

Agenda

- Was ist Deduplizierung?
- Was ist Bacula und was ist Accurate Backup?
- Was sind Base Jobs?
- Base Jobs in der Praxis.

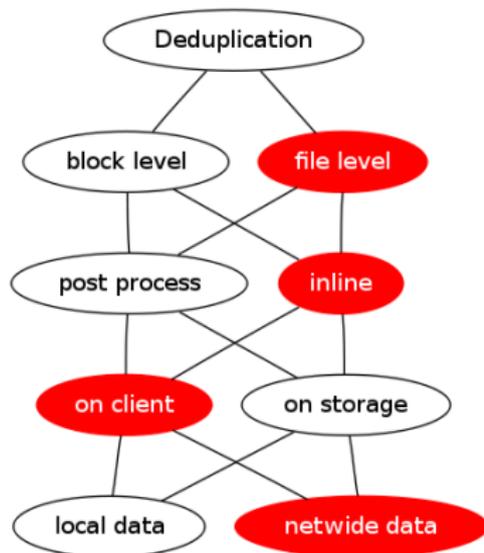
Klassifizierung der Deduplizierungstechniken

- Wie wird dedupliziert?
- Wann wird dedupliziert?
- Wo wird dedupliziert?
- Was wird dedupliziert?



Klassifizierung Bacula Base Jobs

- auf Dateiebene
- während der Sicherung
- auf dem Client
- Daten werden netzweit dedupliziert



- Open-Source (GPL)
- Projektstart: 2000
- netzwerkbasiert
- skalierbar
- unterstützt alle gängigen Betriebssysteme
- Code Qualität und Dokumentation sind hervorragend

Aufgaben eines Backup Systems

- Lesen und schreiben auf dem Endgerät
- Lesen und schreiben auf Sicherungsmedien
- Sicherungen planen und katalogisieren
- (Rück-)Sicherungen durchzuführen über Anwenderschnittstelle

Daten Lesen und schreiben: File Daemon

- Lesen, Schreiben und Verifizieren von Dateien
- Lesen und Schreiben von ACLs, Attributen
- Ausführen von VSS Snapshots unter Windows
- Berechnung von Prüfsummen
- ...



Windows File Daemon



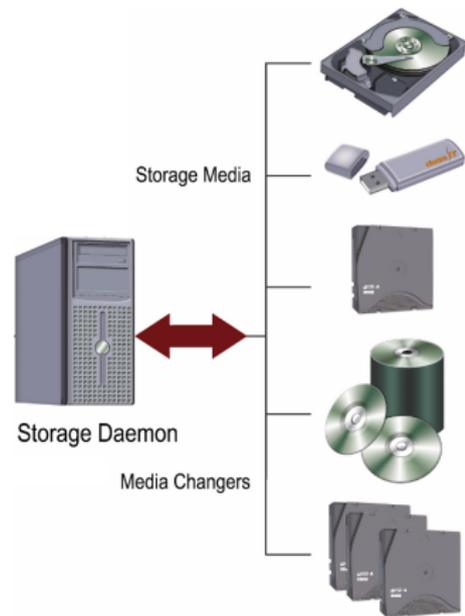
Unix File Daemon



Mac File Daemon

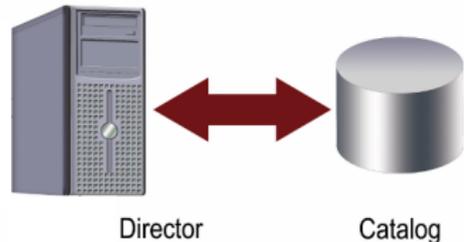
Lesen und schreiben auf Sicherungsmedien: Storage Daemon

- Ansteuerung von Laufwerken: Disk, Tape, WORM, DVD, etc.
- Ansteuerung von Medienwechslern
- Lesen von physischen Medienlabels (Barcodes)
- Schreiben von logischen Bacula-Labels auf die Medien
- Durchführung von Kopien, Migrationen und virtuellen Backups
- Medienfehler verarbeiten



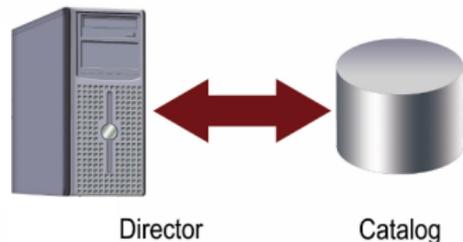
Sicherungen planen und katalogisieren: Bacula Director

- Verwaltung des Bacula Catalog
- Medienverwaltung, Pools
- Zeitplanung
- Konfiguration der zu sichernden Dateien
- Konfiguration der Backupstrategie: voll, inkrementell, differenziell
- Meldungen, Statistiken und Berichte
- Ausführung von Skripten



Bacula Catalog

- Speichert Informationen über alle Dateien, Medien und Jobs
- Backend: MySQL, PostgreSQL oder SQLite
- Umfangreiche Suchmöglichkeit vorhanden
- Über sql beliebige Suchen möglich



Bacula Console

- Anwenderschnittstelle für Rücksicherungen
- Anzeige des Zustands des Systems
- Anzeige des Catalogs
- Starten von Sicherungen



Sitzung Bearbeiten Ansicht Lesezeichen Einstellungen Hilfe

```
BaculaKonf:~ # bconsole
Connecting to Director localhost:9101
bconsole: bsock.c:228 Socket open error. proto=10 port=9101. ERR=Die Adressfamilie
ird von der Protokollfamilie nicht unterstützt
1000 OK: bacula-dir Version: 3.0.2 (18 July 2009)
Enter a period to cancel a command.
*status dir
bacula-dir Version: 3.0.2 (18 July 2009) x86_64-suse-linux-gnu suse 5.x
Daemon started 16-Sep-09 15:03, 0 Jobs run since started.
Heap: heap=278,528 sbytes=68,205 max_bytes=68,782 bufs=243 max_bufs=262

Scheduled Jobs:
Level      Type      Pri  Scheduled      Name      Volume
-----
Incremental Backup 10 16-Sep-09 23:05 BackupClient1 *unknown*
Full      Backup 10 16-Sep-09 23:10 BackupClientErol A80017L4
Full      Backup 11 16-Sep-09 23:10 BackupCatalog *unknown*
****

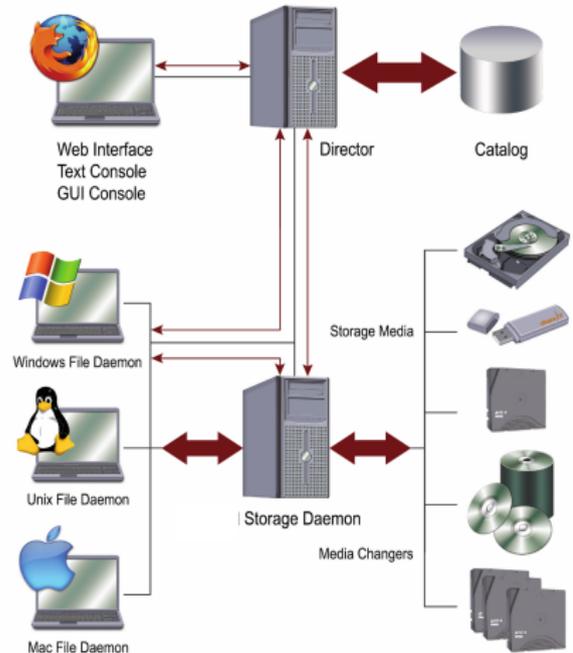
Running Jobs:
Console connected at 16-Sep-09 15:08
No Jobs running.
****

Terminated Jobs:
JobId  Level  Files  Bytes  Status  Finished  Name
-----
5      1      28 70 K  OK      09-Sep-09 14:15 RestoreFiles
6      Full  1      85.42 K  OK      09-Sep-09 14:16 BackupCatalog
9      Full  0      0      Error   10-Sep-09 15:10 BackupClientErol
10     Full  0      0      Error   10-Sep-09 15:30 BackupClientErol
11     Full  0      0      Error   10-Sep-09 15:32 BackupClientErol
12     Full  495   31.17 M  OK      10-Sep-09 15:37 BackupClientErol
13     Full  495   31.17 M  OK      10-Sep-09 15:39 BackupClientErol
14     Full  62,261 4.014 G  Error   10-Sep-09 17:25 BackupClientErol
15     Full  9,351 4.084 G  OK      10-Sep-09 17:35 BackupClientErol
16     Incr  0      0      Error   16-Sep-09 15:03 BackupClient1
****
```

Shell Shell Nr. 2 Shell Nr. 3

Das Gesamtbild

- Kommunikation erfolgt über TCP/IP
- Es werden definierte Ports genutzt
- Kommunikation kann per TLS gesichert werden.



- Vollbackup
 - großer zeitliche Abstand, z.B. Monat
 - Alles wird gesichert
- Differenzielle Sicherung
 - mittlerer zeitlicher Abstand, z.B. Woche
 - Differenz zur letzten Vollsicherung
- Inkrementelle Sicherung
 - kleiner zeitlicher Abstand, z.B. Tag
 - Differenz zur letzten beliebigen Sicherung

Wie wird die Differenz festgestellt?

- Jede Datei besitzt Zeitstempelinformationen
- Zeitstempel wird bei Veränderung automatisch neu gesetzt
- Referenzsicherung hat einen bestimmten Zeitstempel
- Differenz = allen Dateien, deren Zeitstempel neuer als der Referenzzeitpunkt ist.

Zusammenfassung Zeitstempelbasierte Sicherungstechnik

- Differenziellen / Inkrementelle Sicherung:
 - Es gibt einen Referenzzeitpunkt.
 - Dies ist der Zeitpunkt der vorherige Sicherung.
 - Es werden alle Dateien gesichert, bei denen der Zeitstempel neuer ist als der Referenzzeitpunkt.

Ursache

- Datei wird mit altem Zeitstempel erzeugt
- Dateibaum wird verschoben
- Programm manipuliert Zeitstempel (z.B. Virens Scanner)

Wirkung

- Zu sichernde Dateien werden nicht gesichert
- Nicht zu sichernde Dateien werden gesichert
- Gelöschte Dateien werden wiederhergestellt

Lösung: Accurate Backup

- Übertragung der Informationen zu bekannten Dateien und Verzeichnissen an den Client
- Zeitstempelprobleme vermeiden durch Vergleich von:
 - Inodes
 - Zugriffsrechte
 - Anzahl der Links
 - Benutzer und GruppenID
 - Größe
 - Zugriffszeit
 - Modificationszeit Änderungszeit
 - MD5 und SHA1 Signatur.
- Nachteil: Ressourcenintensiv

Ink/Diff

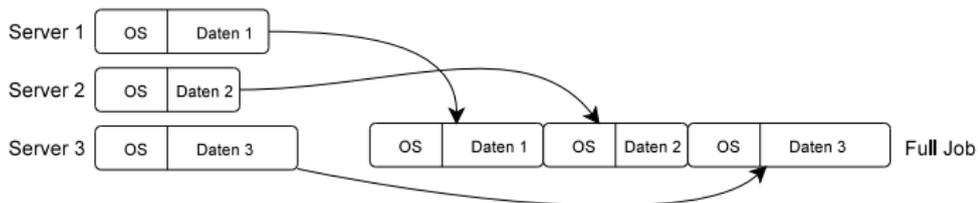
- Alle Dateien werden weggelassen, die bei der Referenzsicherung *dieses Systems* bereits gesichert wurden.

Base Jobs

- Alle Dateien werden weggelassen, die bei der Referenzsicherung *irgendeines Systems* bereits gesichert wurden.

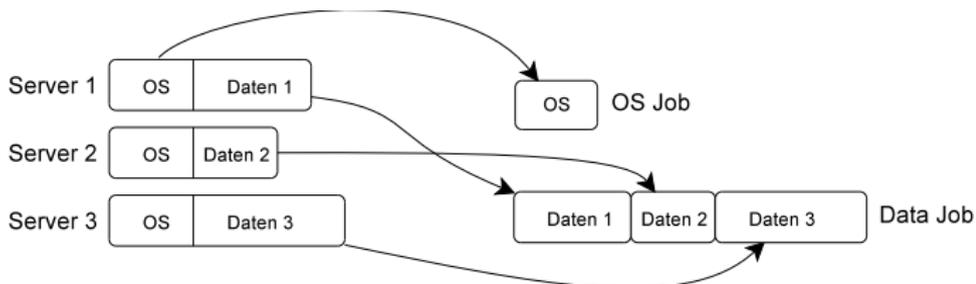
Sicherung identischer Systeme

Bei identischen Systemen werden die Daten und OS mehrfach gesichert.



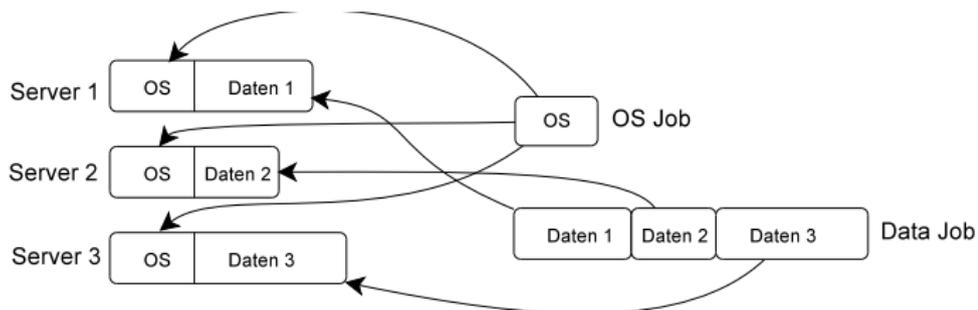
Grundüberlegung bei Basejobs

Bei identischen Systemen wird der gemeinsame Teil nur einmal gesichert.



Rücksicherung bei Basejobs

Bei Rücksicherungen wird der individuelle Teil aus der eigenen Sicherung, der gemeinsame Teil aber aus dem Base Jobs zurückgespielt



Ablauf beim Einsatz von Base Jobs

- 1 Sicherung eines **Base Jobs** in dem möglichst die identischen Daten aller Server gesichert werden.
- 2 "*Normale*" Sicherungen, bei der zu referenzierende Base Job eingestellt wird. *Accurate* ist notwendig.
- 3 Wenn die gemeinsame Basis sich erheblich ändert, erneut einen Base Job ausführen.

- VM1: frosccon-server
 - director daemon
 - storage daemon
 - file daemon
- VM2: frosccon-client
 - file daemon
- Jobs
 - frosccon-base-job
 - frosccon-client-job (basierend auf frosccon-base-job)
 - frosccon-server-job (basierend auf frosccon-base-job)

- 1 Mit dem Einsatz von **Base Jobs** kann das Backupvolumen bei identischen Systemen erheblich verringert werden.
- 2 Bacula Base Jobs arbeiten auf Dateiebene. Für den Einsatz von Base Jobs ist *Accurate Backup* notwendig.
- 3 Gemeinsame Dateien werden nur einmal gespeichert.
- 4 Die Deduplizierung erfolgt bereits auf dem Client.

- 1 Diplomanden gesucht im Bereich Open Source / Linux / C / Python
- 2 Bacula Buch
- 3 Bacula Konferenz

<http://www.dass-it.de>
Philipp.Storz@dass-it.de

dass IT



Fragen?

dass IT

dass IT