

ZFS

Siste ord innen filsystemer

Trond Endrestøl

Fagskolen Innlandet, IT-avdelingen

2. januar 2014

Foredragets filer I

- Filene til foredraget er tilgjengelig gjennom:
 - Subversion: [svn co svn://svn.ximalas.info/zfs-foredrag](svn://svn.ximalas.info/zfs-foredrag)
 - Web: svnweb.ximalas.info/zfs-foredrag
 - Begge metodene er tilgjengelig med både IPv4 og IPv6
- [zfs-foredrag.pdf](#) vises på lerretet
- [zfs-foredrag.handout.pdf](#) er mye bedre for publikum å se på
- [zfs-foredrag.handout.2on1.pdf](#) og [zfs-foredrag.handout.4on1.pdf](#) er begge velegnet til utskrift
- *.169.pdf-filene er i 16:9-format
- *.1610.pdf-filene er i 16:10-format

Foredragets filer II

- Foredraget er mekket ved hjelp av [GNU Emacs](#), [AUCTEX](#), [pdfTEX](#) fra [MiKTEX](#), [LATEX](#)-dokumentklassa [beamer](#), [Subversion](#), [TortoiseSVN](#) og [Adobe Reader](#)
- Hovedfila bærer denne identifikasjonen:
\$Ximalas: trunk/zfs-foredrag.tex 8 2014-01-02 11:53:00Z trond \$
- Driverfila for denne PDF-fila bærer denne identifikasjonen:
\$Ximalas: trunk/zfs-foredrag.handout.4on1.169.tex 3 2013-12-23 13:42:53Z
trond \$
- Copyright © 2013 Trond Endrestøl
- Dette verket er lisensiert med: [Creative Commons, Navngivelse-DelPåSammeVilkår 3.0 Norge \(CC BY-SA 3.0\)](#) 

Oversikt over hele foredraget

Del 1: ZFS?

- 1 Hva er ZFS?
- 2 Hva er grensene til ZFS?
- 3 Hvordan virker ZFS?
- 4 ZFS og RAID-kontrollere
- 5 Hvor kommer ZFS fra?
- 6 Versjonsnummer i ZFS
 - Pool-versjonsnummer
 - Filsystem-versjonsnummer
- 7 Fremtiden for ZFS?

Oversikt over hele foredraget

Del 2: ZFS!

8 Administrasjon av ZFS

- zpool
- zfs

9 Oppretting av pooler

- Enkle pool-eksempler
- Avanserte pool-eksempler

10 zpool-egenskaper

11 zfs-egenskaper

Del I

ZFS?

Oversikt over del 1: ZFS?

1 Hva er ZFS?

2 Hva er grensene til ZFS?

3 Hvordan virker ZFS?

4 ZFS og RAID-kontrollere

5 Hvor kommer ZFS fra?

6 Versjonsnummer i ZFS

- Pool-versjonsnummer
- Filsystem-versjonsnummer

7 Fremtiden for ZFS?

Hva er ZFS?

• ZFS er

- ① Logisk volumhåndterer (Logical Volume Manager, LVM)
- ② Filsystem med bl.a. snapshots, kloner, kompresjon og deduplisering
- ③ Tilbyr også «zvolumer» som lagringsenheter for andre filsystemer

• ZFS tar dataintegritet på alvor; hastighet kommer i senere rekker

• Enklere organisering enn «Storage Spaces» i Microsoft Windows Server 2012

• Lagringen organiseres i pooler som kan bestå av

- ① Enkeldisker/partisjoner
- ② Striping (RAID 0) mellom to eller flere disker/partisjoner
- ③ Speiling (RAID 1) mellom to eller flere disker/partisjoner
- ④ raidz1 (RAID 5, enkel paritet) over tre eller flere disker/partisjoner
- ⑤ raidz2 (RAID 6, dobbel paritet) over fire eller flere disker/partisjoner
- ⑥ raidz3 («RAID 7», trippel paritet) over fem eller flere disker/partisjoner

• Visse kombinasjoner av det overstående er også mulig

Hva er grensene til ZFS?

- ZFS er stort sett grenseløs
 - 128-bit diskadresser
 - Maks. 2^{48} poster i hver katalog
 - Maks. 2^{64} bytes (16 EiB, 16 exabytes) for hver fil
 - Maks. 2^{64} bytes for hvert attributt
 - Maks. 2^{78} bytes (256 ZiB, 256 zebabytes) i hver pool
 - Maks. 2^{56} attributter per fil (egentlig begrenset til 2^{48} attributter)
 - Maks. 2^{64} enheter tilknyttet en gitt pool
 - Maks. 2^{64} pooler i et og samme system
 - Maks. 2^{64} filsystemer i samme pool
 - Ref.: <http://en.wikipedia.org/wiki/ZFS>
- Vis meg det systemet som klarer å sprengne noen av disse grensene!

Hvordan virker ZFS?

- ZFS unngår RAID 5-skrivehullet til eldre RAID-kontrollere som
 - ① Skriver nye data til de samme datablokkene som tidligere
 - ② Leser gamle, urørte data fra de samme datablokkene
 - ③ Regner ut ny paritet for datablokkene
 - ④ Skriver oppdatert paritet til de samme paritetsblokkene som tidligere
 - Hva skjer nå og senere hvis du får strømbrudd mellom punktene 1 og 4?
 - Har diskkontrolleren batteribeskyttet minne?
- ZFS skriver fulle stripere; data og paritet samtidig
- ZFS bruker «copy-on-write»; skriver nye data til ledige diskblokker
- Endringer som hører sammen, samles i transaksjonsgrupper
- Sjekksummer brukes for alt som blir lagret
 - ZFS kontrollerer at leste data er de samme som ble skrevet
 - Oppdages avvik, leter ZFS etter alternativer
 - Finnes alternativer, enten speilkopier eller paritet, så
 - ① Leveres korrekte data til applikasjonen, og
 - ② Avviket korrigeres automatisk på den syke diskens
 - Finnes ingen alternativer, så må filene restaureres fra backup

ZFS og RAID-kontrollere

- Ikke bruk ZFS sammen med RAID-kontrollere!
- RAID-kontrolleren kan i verste fall motarbeide ZFS
- Sett RAID-kontrolleren i JBOD-modus, eller
- La hver harddisk være sitt enslige RAID 0-volum

Hvor kommer ZFS fra?

- Utviklet av Jeffrey Bonwick, Matthew Ahrens og flere kollegaer ved Sun Microsystems, Inc.
- Arbeidet begynte i 2001
- Første prototyp ble ferdig 31. oktober 2001 (halloween)
- ZFS → Solaris, oktober 2005
- ZFS er lisensiert etter «Common Development and Distribution License» (CDDL)
- ZFS → OpenSolaris, november 2005
- ZFS → FreeBSD, april 2007
- Linux' GPL v2-lisens kompliserer import av ZFS
 - ZFS i Linux gjennom FUSE gjenstår som en (treg) mulighet
 - Brian Behlendorf ved Lawrence Livermore National Laboratory (LLNL) har laget «Native ZFS for/on Linux»
- ZFS var tilgjengelig i Mac OS X 10.5, bare read-only, men har vært tilbaketrukket siden oktober 2009
- Noen Mac OS X-entusiaster har laget sine egne ZFS-varianter
- Andre OS med ZFS-støtte: OpenIndiana, FreeNAS, PC-BSD, GNU/kFreeBSD og NetBSD

Versjonsnummer i ZFS

- Pool-versjonene 1–28 og filsystem-versjonene 1–5 er tilgjengelig gjennom OpenSolaris og illumos
- Pool-versjonene 29–34 og filsystem-versjon 6 er bare tilgjengelig i Solaris 11 (Express)
- OpenSolaris har gått videre til feature-flags og pool-versjon 1000
- illumos har gått videre til feature-flags og pool-versjon 5000
- De fleste OS-er utenom Solaris, samarbeider om videreutviklingen av illumos-varianten

Versjonsnummer i ZFS I

Pool-versjonsnummer

- ❶ First release
- ❷ Ditto Blocks
- ❸ Hot spares, double-parity RAID-Z (raidz2), improved RAID-Z accounting
- ❹ zpool history
- ❺ gzip compression for ZFS datasets
- ❻ "bootfspool property
- ❾ ZIL: adds the capability to specify a separate Intent Log device or devices
- ❿ ability to delegate `zfs(1M)` administrative tasks to ordinary users
- ❾ CIFS server support, dataset quotas
- ❽ Devices can be added to a storage pool as "cache devices"
- ❾ Improved zpool scrub/resilver performance

Versjonsnummer i ZFS II

Pool-versjonsnummer

- ❿ Snapshot properties
- ❾ Properties: `usedbysnapshots`, `usedbychildren`, `usedbyreservation`, and `usedbydataset`
- ❿ passthrough-x aclinherit property support
- ❾ Properties: `userquota`, `groupquota`, `userused` and `groupused`; also required FS v4
- ❿ STMF property support
- ❿ triple-parity RAID-Z
- ❿ ZFS snapshot holds
- ❿ ZFS log device removal
- ❿ zle compression algorithm that is needed to support the ZFS deduplication properties in ZFS pool version 21, which were released concurrently
- ❿ Deduplication

Versjonsnummer i ZFS III

Pool-versjonsnummer

- ❿ `zfs receive properties`
- ❿ slim ZIL
- ❿ System attributes. Symlinks now their own object type. Also requires FS v5.
- ❿ Improved pool scrubbing and resilvering statistics
- ❿ Improved snapshot deletion performance
- ❿ Improved snapshot creation performance (particularly recursive snapshots)
- ❿ Multiple virtual device replacements
- ❿ RAID-Z/mirror hybrid allocator
- ❿ ZFS encryption
- ❿ Improved 'zfs list' performance
- ❿ One MB block support
- ❿ Improved share support
- ❿ Sharing with inheritance

Versjonsnummer i ZFS I

Filsystem-versjonsnummer

- ❶ First release
- ❷ Enhanced directory entries. In particular, directory entries now store the object type. For example, file, directory, named pipe, and so on, in addition to the object number.
- ❸ Support for sharing ZFS file systems over SMB. Case insensitivity support. System attribute support. Integrated anti-virus support.
- ❹ Properties: userquota, groupquota, userused and groupused
- ❺ System attributes; symlinks now their own object type
- ❻ Multilevel file system support

Fremtiden for ZFS?

- Oracle kjøpte opp Sun Microsystems, Inc., 27. januar 2010
- Oracle gjorde OpenSolaris om til «ClosedSolaris» i mai 2010
- Hele ZFS-teamet hos Oracle sa opp på dagen, omtrent 90 dager etter denne avgjørelsen ifølge Bryan Cantrill
- ZFS lever videre hos
 - Oracle Solaris
 - illumos/OpenZFS
 - OpenIndiana
 - FreeBSD
 - Delphix
 - iXsystems
 - Joyent
 - NetBSD
 - Nexenta
 - Linux

Del II

ZFS!

Oversikt over del 2: ZFS!

❸ Administrasjon av ZFS

- zpool
- zfs

❹ Oppretting av pooler

- Enkle pool-eksempler
- Avanserte pool-eksempler

❽ zpool-egenskaper

❾ zfs-egenskaper

Administrasjon av ZFS

- To kommandoer (med underkommandoer):
 - ① **zpool**
 - Administrasjon av lagringspoolene
 - ② **zfs**
 - Administrasjon av filsystemer, snapshots, kloner, m.m.
- Det finnes en tredje kommando for de nysgjerrige: **zdb**
 - Brukes for å avlese indre ZFS-detaljer

Administrasjon av ZFS I

zpool-kommandoer

- **zpool add**
- **zpool attach**
- **zpool clear**
- **zpool create**
- **zpool destroy**
- **zpool detach**
- **zpool export**
- **zpool get**
- **zpool history**
- **zpool import**
- **zpool iostat**

Administrasjon av ZFS II

zpool-kommandoer

- **zpool labelclear**
- **zpool list**
- **zpool offline**
- **zpool online**
- **zpool reguid**
- **zpool remove**
- **zpool reopen**
- **zpool replace**
- **zpool scrub**
- **zpool set**
- **zpool split**
- **zpool status**
- **zpool upgrade**

Administrasjon av ZFS I

zfs-kommandoer

- **zfs allow**
- **zfs bookmark**
- **zfs clone**
- **zfs create**
- **zfs destroy**
- **zfs diff**
- **zfs get**
- **zfs groupspace**
- **zfs holds**
- **zfs hold**
- **zfs inherit**

Administrasjon av ZFS II

zfs-kommandoer

- zfs jail
- zfs list
- zfs mount
- zfs promote
- zfs receive
- zfs release
- zfs rename
- zfs rollback
- zfs send
- zfs set
- zfs share

Administrasjon av ZFS III

zfs-kommandoer

- zfs snapshot
- zfs unallow
- zfs unjail
- zfs unmount
- zfs unshare
- zfs upgrade
- zfs userspace

Opprettning av pooler

- zpool create [opsjoner] *navn-på-pool* [*organiseringstype*] *ingredienser* [*organiseringstype* *ingredienser*] ...
- Unngå å plassere mer enn 9 enheter i hver vdev
- I stedet for å stripe en pool over 20 harddisker, vurdér å speile to og to harddisker i 10 grupper

Opprettning av pooler

Enkle pool-eksempler

- Singledisk:
 - zpool create rpool da0
- RAID 0 over to disk:
 - zpool create rpool da0 da1
- RAID 1 over to disk:
 - zpool create rpool **mirror** da0 da1
- RAID 5 over tre disk:
 - zpool create rpool **raidz1** da0 da1 da2
- RAID 6 over fire disk:
 - zpool create rpool **raidz2** da0 da1 da2 da3
- «RAID 7» over fem disk:
 - zpool create rpool **raidz3** da0 da1 da2 da3 da4

Opprettning av pooler

Avanserte pool-eksempler

- RAID 1+0 (3 vdevs á 2 disker):
zpool create rpool **mirror** da0 da1 **mirror** da2 da3 **mirror** da4 da5
- RAID 5+0 (2 vdevs á 3 disker):
zpool create rpool **raidz1** da0 da1 da2 **raidz1** da3 da4 da5
- RAID 6+0 (2 vdevs á 4 disker):
zpool create rpool **raidz2** da0 da1 da2 da3 **raidz2** da4 da5 da6 da7
- RAID 1+5+0 (2 vdevs, 2 og 3 disker):
zpool create rpool **mirror** da0 da1 **raidz1** da2 da3 da4

zpool-egenskaper I

- size
- capacity
- altroot
- health
- guid
- version
- bootfs
- delegation
- autoreplace
- cachefile
- failmode
- listsnapshots

zpool-egenskaper II

- autoexpand
- dedupditto
- dedupratio
- free
- allocated
- readonly
- comment
- expandsize
- freeing
- feature@async_destroy
- feature@empty_bpobj
- feature@lz4_compress

zpool-egenskaper III

- feature@multi_vdev_crash_dump
- feature@spacemap_histogram
- feature@extensible_dataset

zfs-egenskaper I

- type
- creation
- used
- available
- referenced
- compressratio
- mounted
- quota
- reservation
- recordsize
- mountpoint
- sharenfs

zfs-egenskaper II

- checksum
- compression
- atime
- devices
- exec
- setuid
- readonly
- jailed
- snapdir
- aclmode
- acinherit
- canmount

zfs-egenskaper III

- xattr
- copies
- version
- utf8only
- normalization
- casesensitivity
- vscan
- nbmand
- sharesmb
- refquota
- refreservation
- primarycache

zfs-egenskaper IV

- secondarycache
- usedbysnapshots
- usedbydataset
- usedbychildren
- usedbyreservation
- logbias
- dedup
- mlslabel
- sync
- refcompressratio
- written
- logicalused
- logicalreferenced