

## ZFS - Das Dateisystem der nächsten Generation?



Secure Linux  
Administration Conference  
12.12.2008 Magdeburg

Markus Schreier  
ORDIX AG, Wiesbaden  
ms@ordix.de  
www.ordix.de

# ZFS - Das Dateisystem der nächsten Generation?

Wissen, was zu tun ist.

3. Secure Linux Administration Conference

+ Sonntag: „Rechtsprobleme für Administratoren“

10. bis 12. Dezember 2008



# Agenda



- Überblick
- Ziel: Größenbegrenzungen aufheben
- Ziel: Datenintegrität
- Ziel: Administrierbarkeit
- Ziel: neue Funktionalität
- Ziel: Geschwindigkeit
- Implementierung unter Linux
- Fazit

# Was ist ZFS?



- Ein Dateisystem
- Ein Logical Volume Manager
- Von SUN entwickelt
- Frei unter der CDDL  
(Common Development and Distribution License)
- Verfügbar auf
  - Solaris 10
  - FreeBSD
  - MacOS X <sup>1)</sup>
  - Linux <sup>2)</sup>

1) nur beta

2) mit FUSE

ZFS ist keine Erweiterung, sondern eine völlige Neuentwicklung

- Erstes 128-Bit-Filesystem
- Storage Pool basierend
- Unterstützt Kompression, Quota, Platzreservierung, ACLs, NFSv4
- Gewährleistung der Konsistenz („copy on write“)
- Transaktions-orientiert - immer konsistent
- Einfache Administration (...)
- GUI (Java Web Console)

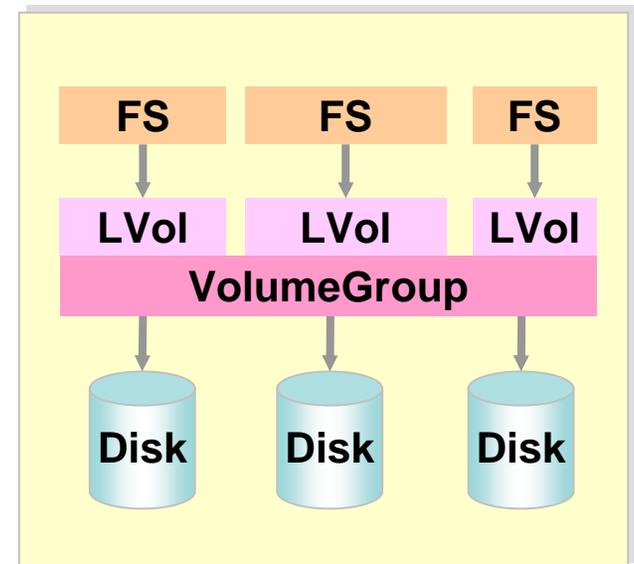
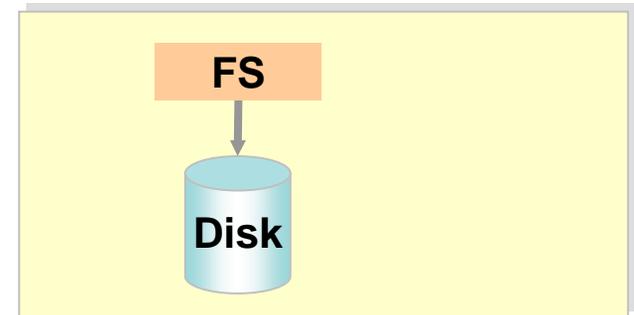
# Linux Filesysteme

Name	verfügbar seit	Hersteller/Anbieter	Typ
ext2	1993	Linux	herkömmlich
ext3	1999	Linux	journaling
ext4	2006	Linux	journaling
reiserFS	2001	Linux	journaling
OCFS /OCFS2	2002	Oracle Linux	cluster-FS
XFS	1994	SGI irix	journaling
JFS	1991	IBM aix	journaling
VxFS	1991	Veritas	journaling
btrfs	Entwicklungsbeginn 2008	Oracle Linux	COW, 64bit



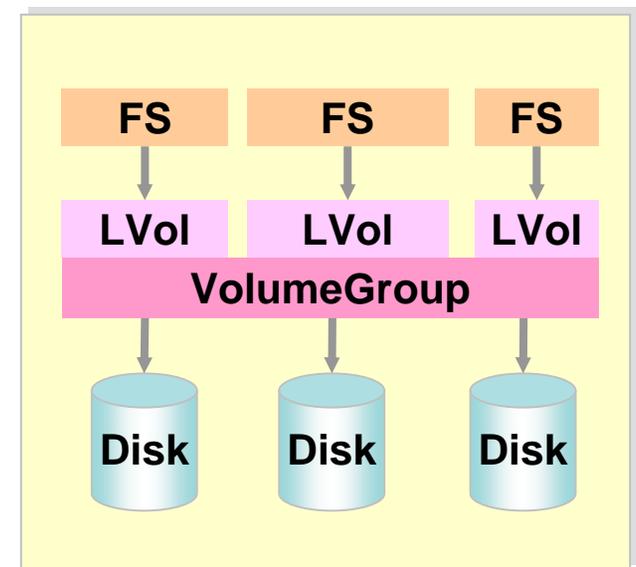
btrfs ist die einzige ernstzunehmende Konkurrenz für zfs, jedoch noch in einem sehr frühen Entwicklungsstadium

- Dateisysteme (z. B. ext2) direkt auf Platten/Devices
- Hinzu kamen neue Ziele:
  - Inkonsistenzen nach System-Crash  
Æ Journaling
  - Mehr Plattenplatz  
Æ striping, concatenating
  - Mehr Sicherheit  
Æ mirror, raid5
  - Mehr Flexibilität  
Æ Softpartitions, LVM, resizefs
  - Mehr Komfort  
Æ LVM (Snapshots, Clone, Komprimierung)



## Nachteile:

- journal kuriert nur das Symptom
- Filesystem kennt Volume-Struktur nicht
- IO-Sequenzen müssen Reihenfolge beibehalten
- Administration schwierig, mehrstufig (fdisk, lvm, fs)
- Keine Erkennung von korrupten Daten



# Ziele von ZFS



- Größenbegrenzungen aufheben
  - erstes 128-Bit-Dateisystem
- Datenintegrität
  - erkennt “silent data corruption”
  - Transaktions-orientiert
  - immer konsistent
  - selbstheilend
- Administrierbarkeit
  - integrierter Logical Volume Manager
- Neue Funktionalität
  - Snapshots, Clone
  - ACLs, xAttr, NFS, ...
- Geschwindigkeit

## Wie erreicht?

- Ursache der bisherigen Probleme erkannt
- Historie über Board geworfen
- Von bestehenden Annahmen gelöst
- Komplette Neuentwicklung

- Überblick
- Ziel: Größenbegrenzungen aufheben
- Ziel: Datenintegrität
- Ziel: Administrierbarkeit
- Ziel: neue Funktionalität
- Ziel: Geschwindigkeit
- Implementierung unter Linux
- Fazit

# Neue Größenordnungen



- In 10 - 15 Jahren reicht 64-Bit-Adressraum nicht aus
  - Æ 128-Bit-FS
  - Æ maximale Größe: 16 ExaByte
- Dynamische Allokation von Metadaten
  - Æ skalierbar

Wert	10er Potenz	Präfix
$2^0 = 1$	-	
$2^{10} = 1.024$	$\sim 10^3$	k (kilo)
$2^{20} = 1.048.576$	$\sim 10^6$	M (Mega)
$2^{30} = 1.073.741.824$	$\sim 10^9$	G (Giga)
$2^{40} = 1.099.511.627.776$	$\sim 10^{12}$	T (Tera)
$2^{50} = 1.125.899.906.842.624$	$\sim 10^{15}$	P (Peta)
$2^{60} = 1.152.921.504.606.846.976$	$\sim 10^{18}$	E (Exa)
$2^{70} = 1.180.591.620.717.411.303.424$	$\sim 10^{21}$	Z (Zetta)
$2^{80} =$	$\sim 10^{24}$	Y (Yotta)

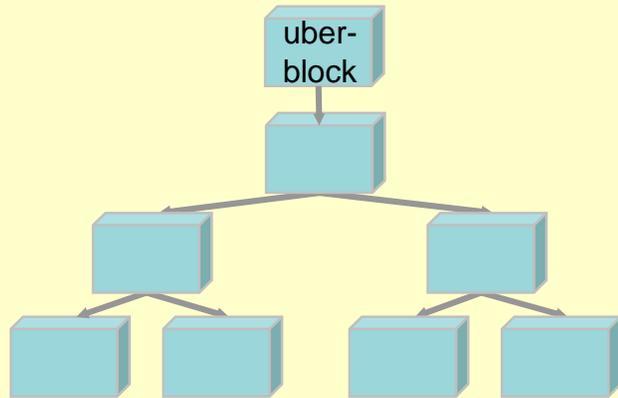
- Überblick
- Ziel: Größenbegrenzungen aufheben
- Ziel: Datenintegrität
- Ziel: Administrierbarkeit
- Ziel: neue Funktionalität
- Ziel: Geschwindigkeit
- Implementierung unter Linux
- Fazit

## Datenintegrität

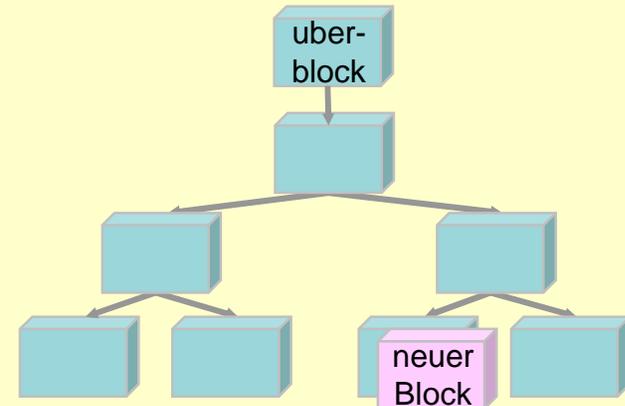
- erkennt “silent data corruption”
- Transaktions-orientiert
- immer konsistent
- selbstheilend

# Konsistenz durch "copy on write"

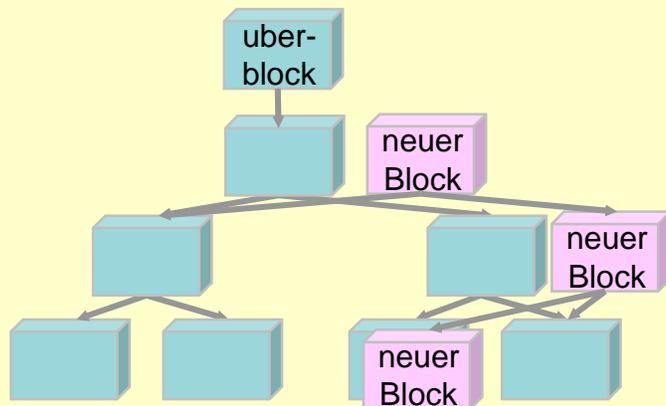
## 1 Ursprünglicher Block-Baum



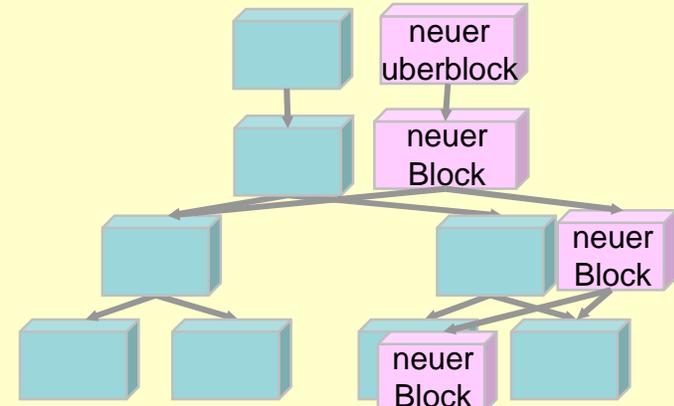
## 2 copy on write eines Datenblocks



## 3 copy on write aller referenzierenden Blöcke



## 4 copy on write des uberblocks (atomar)

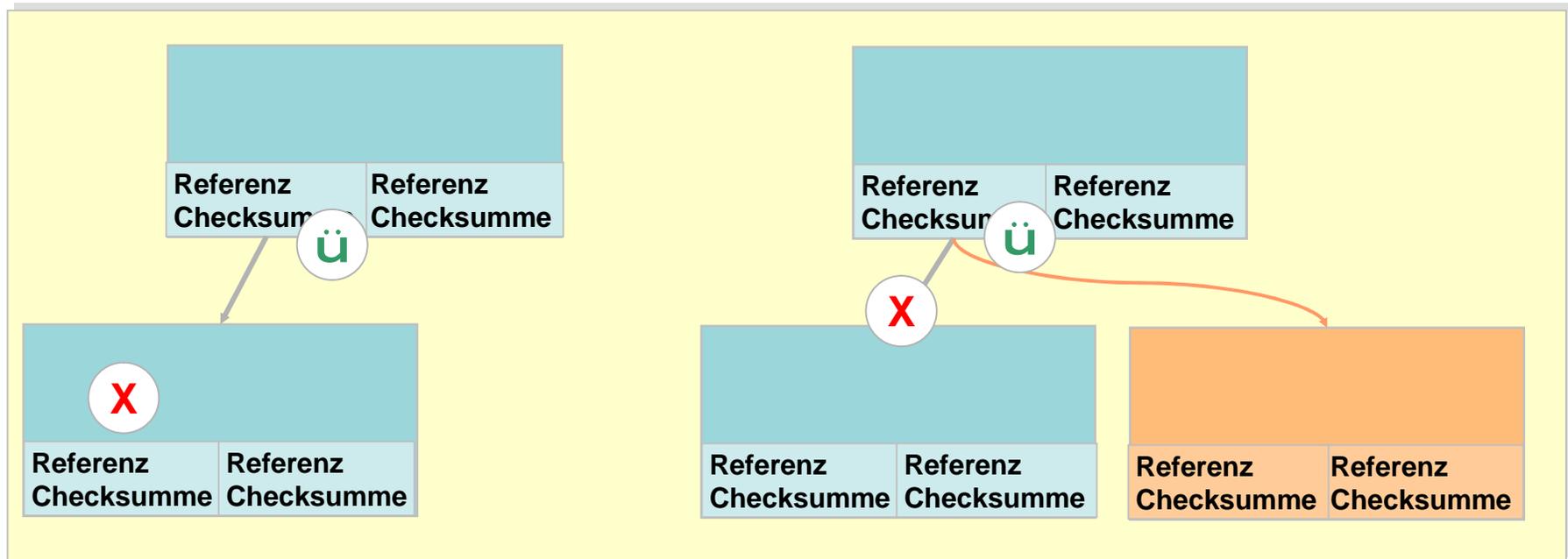


# Konsistenz - Checksummen

- Ende zu Ende Konsistenz
- Checksummen im referenzierenden Block

Æ Fehlerhafte Blöcke werden erkannt

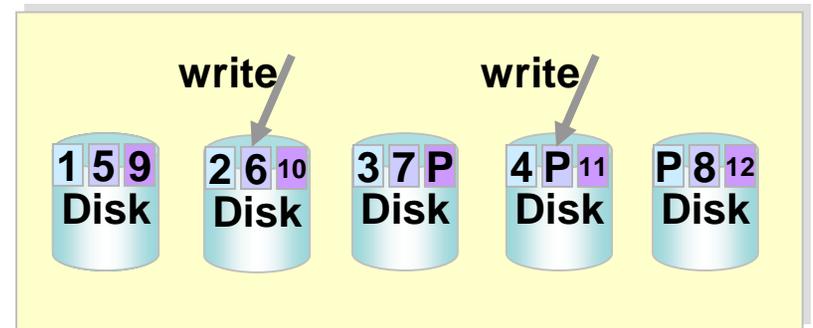
Æ Falsche (aber konsistente) Blöcke werden erkannt



# Konsistenz – RAIDZ / RAIDZ2

## RAID 5 / RAID 6

- Ändern eines Datenblocks
- Ändern von Parity
  - $\text{Æ}$  „RAID5 write hole“



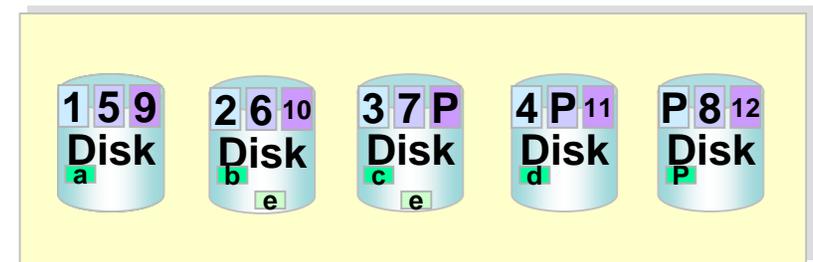
## RAIDZ

- Nutzt copy on write
- Schreiboperationen immer über alle Platten
  - Extendgröße wird angepasst
  - sehr kleine Datenmengen als Spiegel

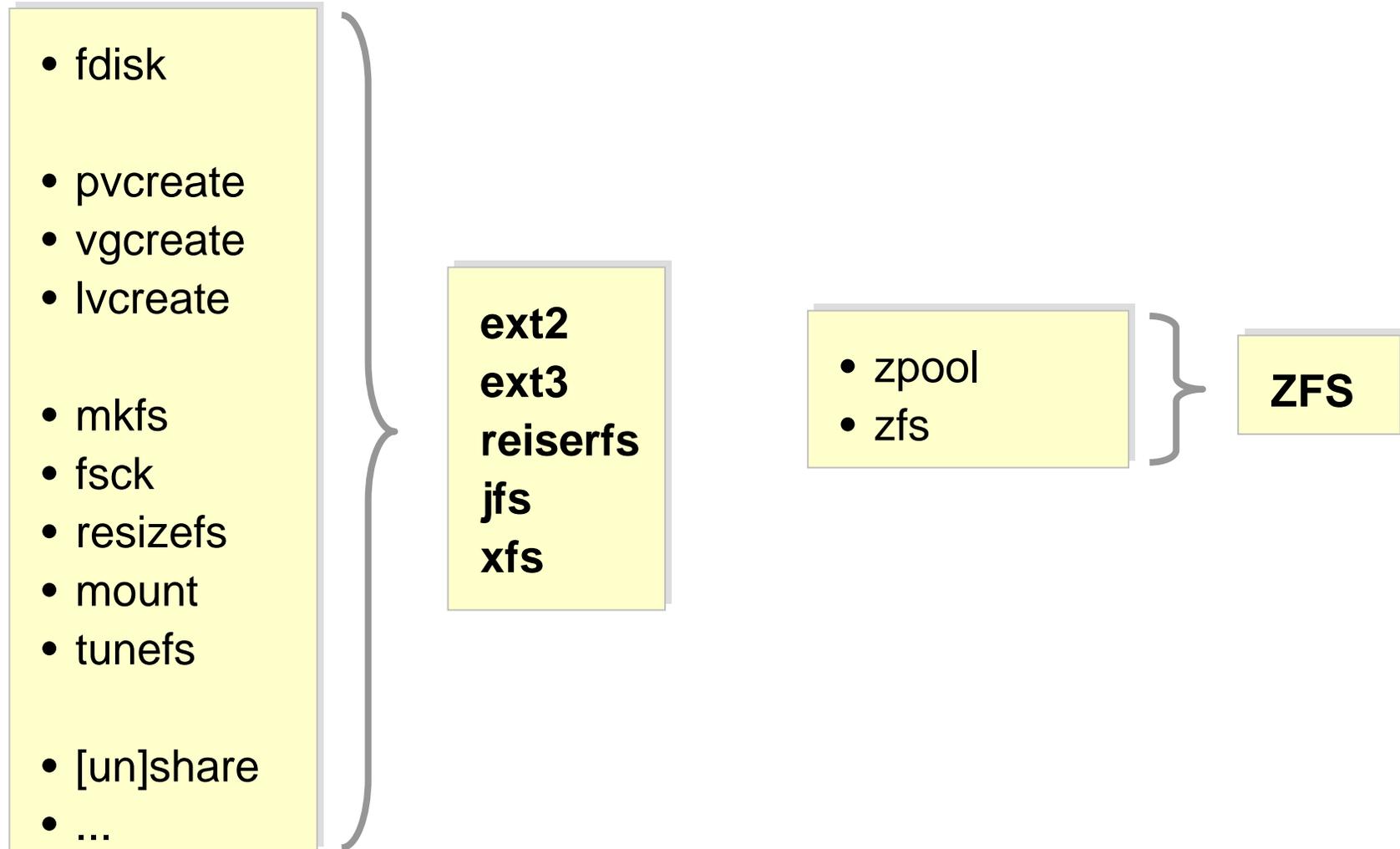
$\text{Æ}$  Keine RAID5 write hole

## RAIDZ2

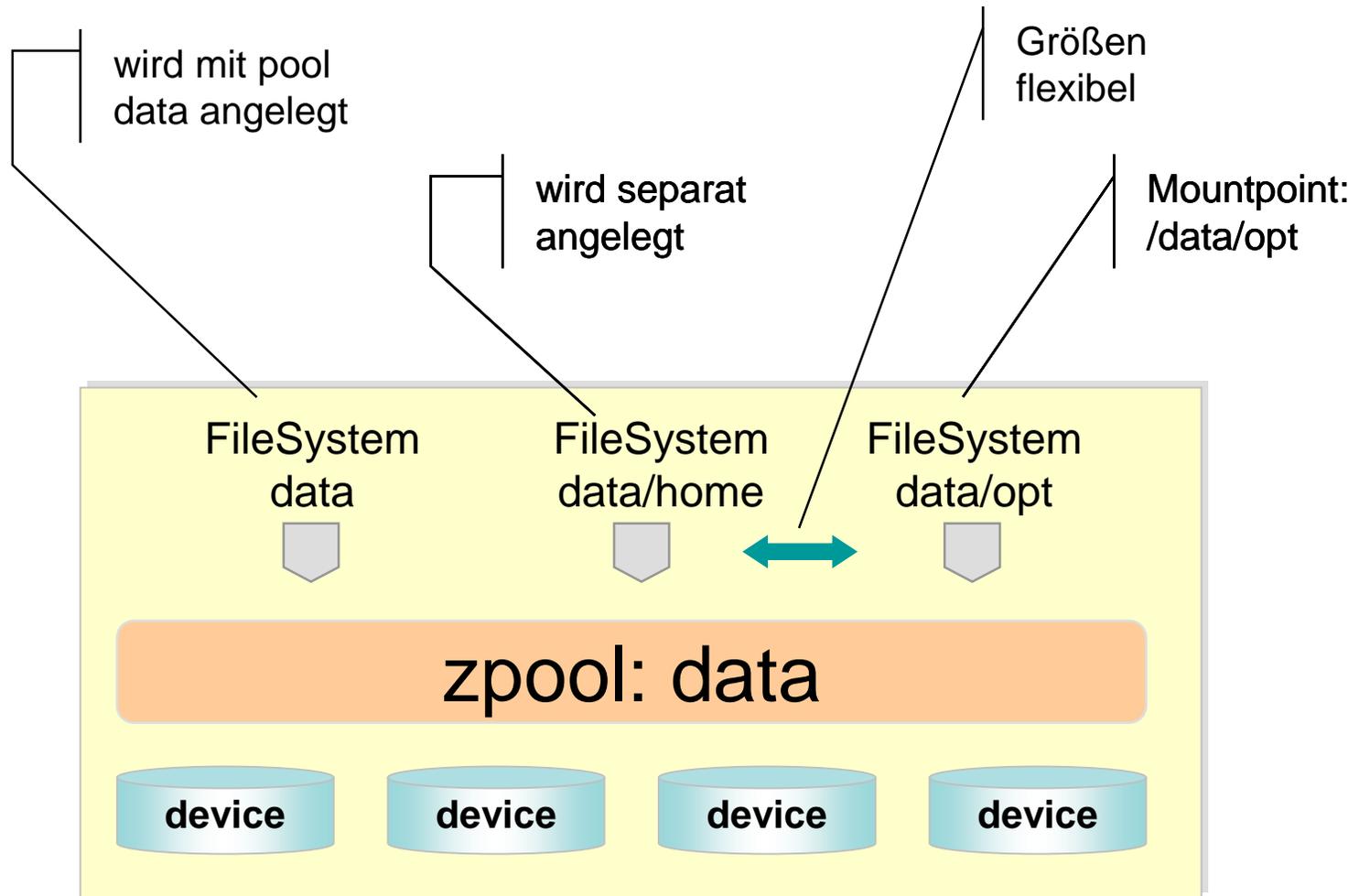
- Wie RAIDZ, doppelte Parity



- Überblick
- Ziel: Größenbegrenzungen aufheben
- Ziel: Datenintegrität
- Ziel: Administrierbarkeit
- Ziel: neue Funktionalität
- Ziel: Geschwindigkeit
- Implementierung unter Linux
- Fazit

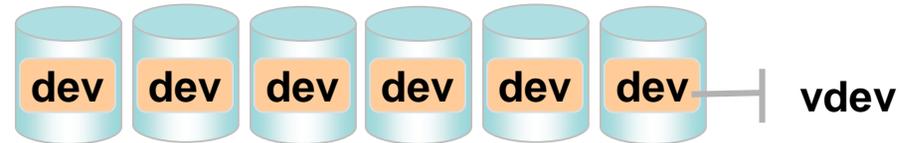


# Storage Pools



# Poolkonfiguration

Stripe/Concat



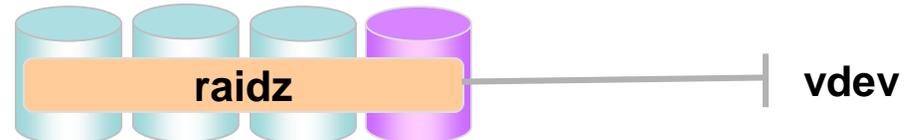
Spiegel



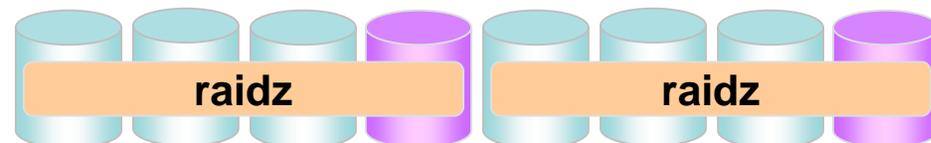
3 Spiegel konkateniert



RaidZ aus 4 Platten



2 konkatenierte RaidZ



RaidZ2 aus 8 Platten



# Pool anlegen



```
zpool create poolname vdev [ vdev[ vdev[ ...]]]
```

*vdev* kann sein

- einfaches Device  
***physdev***
- Spiegel Device  
***mirror physdev physdev [physdev [...]]***
- raidz Device  
***raidz physdev physdev physdev [...]***
- raidz2 Device  
***raidz2 physdev physdev physdev [...]***
- Hotspare für redundante Device  
***spare physdev [physdev [...]]***

```
Bsp.: # zpool create mypool mirror sdb sdc mirror sdd sde
```

# Veränderungen am Plattenlayout



vdev \ Ziel zpool	mehr Redundanz	mehr Platz	weniger Redundanz	weniger Platz
	attach dev	add vdev	dettach dev	delete
ein Device	Spiegel	stripe/concat	-	-
2er-Spiegel	3er-Spiegel	2x 2er-Spiegel	ein Device	-
Xer-Spiegel	(X+1)er-Spiegel	2x Xer-Spiegel	(X-1)er-Spiegel	-
raidz	-	2x raidz	-	-
raidz2	-	2x raidz2	-	-

- Veränderung am Plattenlayout nur eingeschränkt möglich
- Aber: **Wie oft wird das gebraucht?**

- Pool
- Filesystem
  - poolname/fsname
  - poolname/fsname/fsname
- Volume
  - für Block oder Character-Device
  - Name wie bei Filesystem
- Snapshot
  - gehört zu einem Filesystem
    - filesystem@snapshotname

## Beispiele:

- `mypool/home`  
Filesystem
- `mypool/home/ms`  
Filesystem
- `mypool/home@Dienstag`  
Snapshot
- `mypool@gestern`  
Snapshot

- Anlegen eines Pools
  - legt erstes Dateisystem implizit an

```
# zpool create ...
```

- Anlegen eines Dateisystems
  - erstellt mountpoint
  - mounted das Dateisystem
  - kein Eintrag in /etc/fstab nötig

```
# zfs create ...
```

- Steuerung über Dateisystem-Attribute

Reservierung

Komprimierung

Quota

Mountpoint

Beispiel

```
# zfs set reservation=1G mypool/home/ms  
# zfs set quota=1G mypool/home/ms  
# zfs set compression=on mypool/home/ms
```

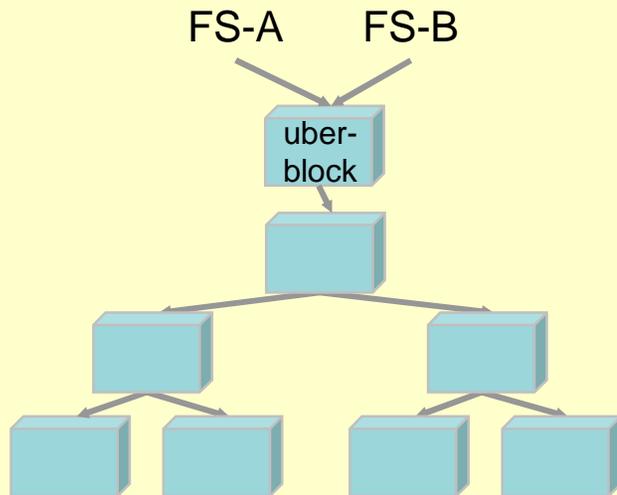
- Überblick
- Ziel: Größenbegrenzungen aufheben
- Ziel: Datenintegrität
- Ziel: Administrierbarkeit
- Ziel: neue Funktionalität
- Ziel: Geschwindigkeit
- Implementierung unter Linux
- Fazit

- Keine Performance-Einbußen
- Anlegen von Snapshots:
  - Atomare Aktion rekursiv über mehrere (Unter)-Filesysteme
- Mittel der Datensicherung
  - `zfs snapshot -r mypool/home@yesterday`
- Lesend Zugreifbar unter <sup>1)</sup>
  - `/mypool/home/.zfs/snapshot/yesterday`
  - `/mypool/home/.zfs/snapshot/2daysago`
- Schreibender Zugriff über Clone

1) noch nicht bei Linux

# Snapshots / Clone & COW

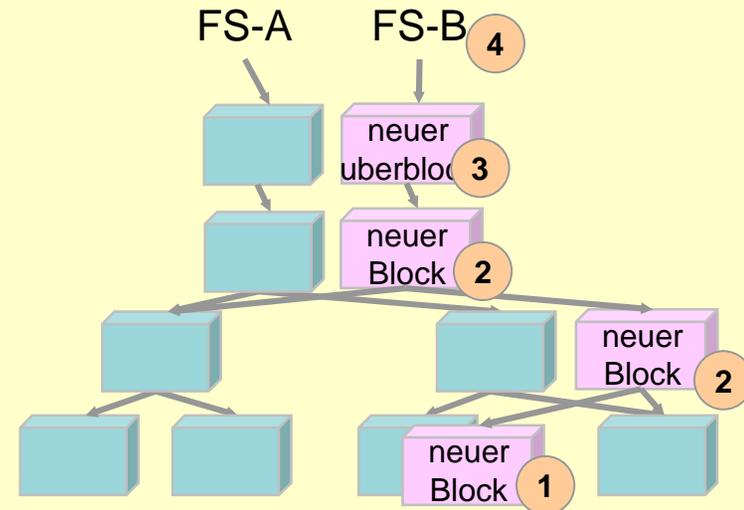
Ursprünglicher Block-Baum



Filesystem A ist Original  
Filesystem B ist Clone

Bisher keine Änderungen  
Alle Filesysteme zeigen  
auf gleichen uberblock

Änderung an FS-B



- 1) im FS-B wird ein Block geändert/neu geschrieben
- 2) darauf referenzierende Blöcke werden neu geschrieben
- 3) uberblock wird neu geschrieben
- 4) Filesystem B verweist auf neuen uberblock  
alle nicht veränderten Blöcke bleiben erhalten

- NFS-Freigaben <sup>1)</sup>
  - Attribut des ZFS-Filesystems
- SMB-Freigaben <sup>2)</sup>
  - Attribut des ZFS-Filesystems
- iSCSI-Targets <sup>1)</sup>
  - benötigt Volumes (/dev/zvol/....) <sup>1)</sup>
  - Attribut des ZFS-Volumes

1) noch nicht bei Linux

2) noch nicht realisiert

- Bisher Posix-ACLs
  - setfacl / getfacl
  - ermöglicht UNIX-Rechte weiteren Benutzern/Gruppen zuzuordnen
- NFSv4-ACLs / Windows -ACLs
  - erweitertes chmod
  - zusätzliche Rechte:  
add\_file (w), add\_subdirectory (p), delete (d),  
delete\_child (D), execute (x), list\_directory (r),  
read\_acl (c), read\_attributes (a), read\_data (r),  
read\_xattr (R), write\_xattr (W), write\_data (w),  
write\_acl (C), write\_attributes (A), write\_owner (o),  
nicht implementiert: synchronize (s), append\_data (p)

# ACL - Vergleich

	<b>posix ACL</b>	<b>Windows, NFSv4 -ACL</b>
<b>Filesysteme</b>	ext2, reiserfs, nfs (bis v3)	ZFS, NFSv4
<b>Kompatibilität</b>	Posix-draft, Linux, ufs	ähnlich NT-Dateirechten
<b>Einstellen</b>	setfacl	chmod
<b>lesen/ansehen</b>	getfacl	ls -v
<b>Umfang</b>	User-Zuordnung erweitert	Rechte differenziert
<b>Transfer</b>	posix -> nfsv4 - ok	nfsv4 -> posix - nur bedingt

- NFSv4 ACLs sind vielfältiger als Posix ACLs
- Besser für Interoperabilität mit Windows
- ACL = Access Control List
- List bedeutet (hier): mehrere Zeilen mit deny/allow-Regeln

- **Erweiterte Attribute (xAttr)**
  - Beliebige Daten/Dateien als Metadaten zu einer Datei ablegen
  - `runat datei [cmd]`
  - wechselt in xattr-Verzeichnis der angegebenen Datei
  - führt dort Kommando `cmd` aus
  - `cp, mv, cat, ls, chmod, chown, ... nutzbar`

- Überblick
- Ziel: Größenbegrenzungen aufheben
- Ziel: Datenintegrität
- Ziel: Administrierbarkeit
- Ziel: neue Funktionalität
- Ziel: Geschwindigkeit
- Implementierung unter Linux
- Fazit

- Transaktionen werden zu Transaktions-Gruppen zusammengefasst
- Intelligente Block-Allozierung (Extents, COW)
- Viele Änderungen ergeben nur eine Schreiboperation
- SUN lieferte viele Performance-Verbesserungen im ZFS-Stack
  
- Nicht innerhalb einer Platte stripen
- Pool über 80 % belegt => Schreib-Performance-Einbußen
- Ganze Platten verwenden IO-Cache der Platte durch ZFS beeinflussbar
  
- Besondere Frage bei FUSE

# Performance - Vergleich



	<b>SLO-Filesystem</b>	<b>COW-Filesystem</b>
<b>Filesystem-Beispiele</b>	statically laid out ext3, reiserfs	copy on write zfs
<b>Random Write in preallocated File</b>	Seek des Blocks schreiben <i>langsam</i>	sequentielles Schreiben <i>schnell</i>
<b>Sequentielles Lesen</b>	Sequentielles IO <i>schnell</i>	Seek des Blocks lesen <i>langsam</i>
<b>Tuning</b>		Große Recordsize => besseres Lesen, schlechteres Schreiben

- Problem für bestimmte Anwendungen
- ZFS in der Entwicklung, Lösungen sind in Arbeit
- Für viele Anwendungen KEIN Problem

- Optimierte Datenbanken  
Oracle, Informix, DB/2, Adabas C, ...
  - Optimiert auf Raw Devices (Kopfbewegungen)
  - **Æ** Daher beste Performance auf Raw Devices
  - **Æ** Rewrite Filesysteme (UFS, VxFS) gehen auch noch
  - ZFS ist eine zusätzliche Schicht (nochmal Log+Writer), die Performance wegnimmt
- Nicht optimierte Datenbanken  
MySQL, PostGreSQL, MaxDB, SQLite
  - ZFS stellt I/O Optimierung für die Datenbank
  - **Æ** Beste Performance auf ZFS

# Agenda



- Überblick
- Ziel: Größenbegrenzungen aufheben
- Ziel: Datenintegrität
- Ziel: Administrierbarkeit
- Ziel: neue Funktionalität
- Ziel: Geschwindigkeit
- Implementierung unter Linux
- Fazit

## Linux

unterliegt  
Gnu Public License



GPL **ist** anerkannt von  
**Open Source Initiative**

## ZFS

unterliegt  
Common Development and  
Distribution License



CDDL **ist** anerkannt von  
**Open Source Initiative**

GPL ist zu CDDL inkompatibel

Æ ZFS ist nicht in den Linux-Kernel integriert

Æ ZFS on FUSE

Ausblick:



Jeff Bonwick und  
Linus Torvalds

# FUSE - Einführung

## Kernelmodul (hier: fuse.ko)

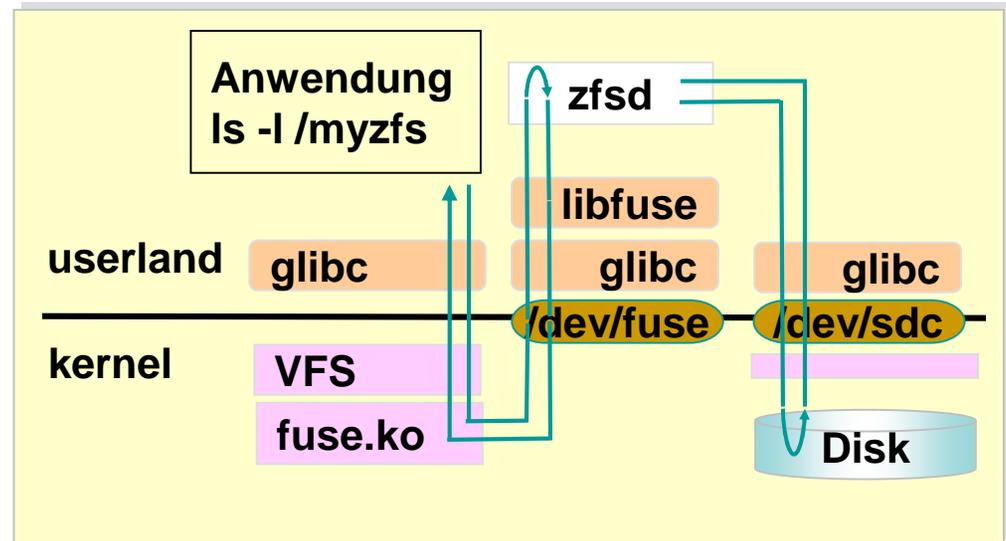
- nutzt Abstraction-Layer VFS ( wie ext3, NFS, ...)
- übersetzt io-Anfragen
- stellt /dev/fuse bereit
- leitet diese an /dev/fuse weiter

## libfuse

- kommuniziert über /dev/fuse
- implementiert calls (open, link, ..)

## zfsd

- user-land –Dämon
- nutzt libfuse
- implementiert eigentliches FS
- greift auf Device zu



- Benötigt Kernelmodul fuse.ko
- Benötigt Dämon zfs-fuse
- Benötigt `zfs mount -a`
  
- Kein Zugriff auf Snapshot über `.zfs` –Verzeichnis
- Rollback von zfs-Snapshot macht kein `umount/mount`  
`Æ` kann zu IO-Fehlern bei der Anwendung kommen
- `promote` von Clonen nimmt nicht alle Snapshots mit
- `zfs set mountpoint=`
  - `mounted` Filesystem nicht wieder automatisch
  - nimmt bestehende Sub-Filesysteme nicht mit
  
- Volumes noch nicht implementiert

# ZFS und Linux - II



sharenfs	spielt mit NFS-Server nicht zusammen
sharesmb	bisher nicht lauffähig (auch bei Solaris)
shareiscsi	benötigt Volumes
acl	statt chmod wird zfs allow/unallow verwendet
xattr	benötigt runat
zfs promote	liefert leicht anderes Ergebnis
FS-Hierarchie	Änderungen an übergeordneten Dateisystem greifen nicht durch

- Überblick
- Ziel: Größenbegrenzungen aufheben
- Ziel: Datenintegrität
- Ziel: Administrierbarkeit
- Ziel: neue Funktionalität
- Ziel: Geschwindigkeit
- Implementierung unter Linux
- Fazit

- YES: - ZFS ist das Dateisystem der nächsten Generation ,denn
  - ZFS Design & Funktionalität sucht seinesgleichen
- NO: ZFS ist in dieser Generation noch nicht produktiv einsetzbar, denn
  - FUSE ist keine produktiv einsetzbare Lösung für ZFS
  - auch unter Solaris muss ZFS bzgl. Performance und Funktion geprüft werden
- Maybe: Es wird vergleichbare Technik unter Linux geben, möglich als
  - ZFS
  - btrfs
  - von Ihnen  geschrieben

Vielen Dank  
für Ihre Aufmerksamkeit!

einfach.gut.beraten.

