

# Verteilte Filesysteme

Ceph, GlusterFS, Client-Sync (DiscoFS), ...

Jochen Schulz, Ralph Krimmel, Max Voit, Piotr Kasprzak

Georg-August Universität Göttingen



## 1 Einführung und Motivation

## 2 Parallel, distributed, fault-tolerant filesystems

- Auswahl verteilte Dateisysteme
- Ceph Theorie
- Ceph Praxis

## 3 Erfahrungen

## 1 Einführung und Motivation

## 2 Parallel, distributed, fault-tolerant filesystems

- Auswahl verteilte Dateisysteme
- Ceph Theorie
- Ceph Praxis

## 3 Erfahrungen

## A: Synchron

- Compute-server: (Fokus: Compute)
- Clients: Laptops mit remote gemounteten Home (Fokus: tägliche Arbeit)
- (File-)Server: (Fokus: shared ressources)

## B: Asynchron

- Daten stets (auch) lokal.
- Asynchroner Datenabgleich mit remote storage (dropbox, powerfolder, git ...)

## Prinzipielle Vor- und Nachteile

- A: Flexibilität im Ort und Gerät (Keine nennenswerte Transferzeit)
- A: Parallelität
- A: Homogenität und Synchrones Arbeiten (z.B. lokales Editieren , remote Ausführen)
- B: Mobilität des Geräts
- B: Daten stets lokal (aber grössere Transferzeit)

# Details bestehendes Setup (Typ A)

- **clients ( $\approx 50$ ):** Linux-Laptops, Datensync per unison
- **Login Server (2):** remote Desktop, ssh, Datentransfer
- **Compute Server (6):** remote Desktop, ssh, Datentransfer, FHGFS, WiRe, CUDA
- **Fileserver (2):** (kein SAN)
  - lokales Dateisystem: ZFS
  - Netzwerkdateisystem: NFS (v3) und samba
  - *Hot spare Replikation* durch AVS

# Details bestehendes Setup (Typ A)

- **clients ( $\approx 50$ ):** Linux-Laptops, Datensync per unison
- **Login Server (2):** remote Desktop, ssh, Datentransfer
- **Compute Server (6):** remote Desktop, ssh, Datentransfer, FHGFS, WiRe, CUDA
- **Fileserver (2):** (kein SAN)
  - lokales Dateisystem: ZFS
  - Netzwerkdateisystem: NFS (v3) und samba
  - *Hot spare Replikation* durch AVS
- **Sonstige Server (6):** Virtualisiert mit OpenVZ
  - Email & Groupware und LDAP
  - Diverse Lizenzserver
  - Nameserver und dhcpd
  - puppet und FAI

# Schwierigkeiten des Setups

- Fileserver
  - ZFS (nur lokal)
  - NFS hat nur einen Zugangsknoten  $\Rightarrow$  ineffiziente Hardware-Nutzung, keine intrinsische Redundanz
  - NFS wegen Sicherheit und des clientssetups nur intern nutzbar
  - **AVS** verursacht Last auf dem Haupt-Server



# Schwierigkeiten des Setups

- Fileserver
  - ZFS (nur lokal)
  - NFS hat nur einen Zugangsknoten  $\Rightarrow$  ineffiziente Hardware-Nutzung, keine intrinsische Redundanz
  - NFS wegen Sicherheit und des clientssetups nur intern nutzbar
  - **AVS** verursacht Last auf dem Haupt-Server
- Sonstige Server
  - Unbenutzte Speicherkapazität
  - Unbenutzte Rechenkapazität (Balancing nur schwer möglich)
  - Unzureichende Verfügbarkeit bei Fehlern (single point of failure)

# Schwierigkeiten des Setups

- Fileserver
  - ZFS (nur lokal)
  - NFS hat nur einen Zugangsknoten  $\Rightarrow$  ineffiziente Hardware-Nutzung, keine intrinsische Redundanz
  - NFS wegen Sicherheit und des clientssetups nur intern nutzbar
  - **AVS** verursacht Last auf dem Haupt-Server
- Sonstige Server
  - Unbenutzte Speicherkapazität
  - Unbenutzte Rechenkapazität (Balancing nur schwer möglich)
  - Unzureichende Verfügbarkeit bei Fehlern (single point of failure)
- Clients
  - Nicht sehr mobil (Mobilität erfordert Restart)
  - Datensynchronisation (Unison) kostet Zeit und Ressourcen bei Login/Logout
  - User muss Vorgang verstehen und beobachten.

Serverseitig: Nutzung von verteilten, fehler-toleranten Dateisystem (Distributed parallel fault-tolerant file systems) und damit storage-cluster und VM-Cluster.

clientseitig:

- Home abkoppeln, nur einzelne Verzeichnisse mounten.
- DiscoFS . FS-Overlay mit vollständigen lokalen Cache (FUSE)
- (geplante) Ceph-erweiterung für disconnected operations

# DiscoFS - DISCOnnected File System (optional)

(Von Robin Martinjak) <https://github.com/rmartinjak/discofs>

## Idee

FS-Overlay über lokalen Daten und remote Daten, die in Sync gehalten werden sollen.

- disconnected operations: Bei disconnect lokal, bei connect lokal und remote.
- Conflict-Handling wenn sich beide Quellen geändert haben.

Komponenten:

- Cache: lokale Daten
- Sync table: Synctime of file objects
- Job: Replikation organisieren
- Database: Job und Sync-Daten speichern
- State check Thread: Online-Status prüfen.
- worker Thread: Jobs abarbeiten

## Idee

Benutze CRUSH um eine zusätzliche Replica eines Homes lokal auf dem Laptop abzulegen.

- Geht die Verbindung verloren sind alle Daten lokal und, theoretisch, kann der lokale Ceph-client weiterarbeiten.
- Das Wiederverbinden ist bereits implementiert, allerdings bräuchte man noch ein Conflict-Handling.

## 1 Einführung und Motivation

## 2 Parallel, distributed, fault-tolerant filesystems

- Auswahl verteilte Dateisysteme
- Ceph Theorie
- Ceph Praxis

## 3 Erfahrungen

# Anforderungen an ein verteiltes Dateisystem

- Hohe Performanz durch Striping
  - Durchsatz
  - Latenz
- Skalierbarkeit
- Verfügbarkeit
- Zuverlässigkeit

**Problem:** Fehler sind mehr die Norm als die Ausnahme

- System günstiger Hardware
- Fehler in Programmen
- Menschliches Versagen
- Betriebssystemfehler
- Ausfall von:
  - Festplatten
  - Arbeitsspeicher
  - Kabel
  - Netzwerk
  - ...



- Objektbasiert
- Festplatten → Intelligent object storage devices (OSD's)
  - CPU
  - Netzwerkschnittstelle
  - Lokaler Cache
- Metadaten getrennt von Nutzdaten

## 1 Einführung und Motivation

## 2 **Parallel, distributed, fault-tolerant filesystems**

- Auswahl verteilte Dateisysteme
  - Ceph Theorie
  - Ceph Praxis

## 3 Erfahrungen

# Auswahl verteilte Dateisysteme

- Lustre** Alter Kernel und nur ein Metadatenserver
- Google FS** Proprietär, ausgelegt auf große Dateien
- GlusterFS** Nur FUSE, keine advanced features
- XtremeFS** nett ;) aber Fokus liegt auf Geo-Features
- FHGS** Proprietärer server, Opensource client. keine Redundanz, im HPC-Umfeld aber gut (Bei uns seit kurzem im Einsatz).
- pNFS** Vielversprechend, bisher aber keine Tests durchgeführt (zu spät?)

## 1 Einführung und Motivation

## 2 **Parallel, distributed, fault-tolerant filesystems**

- Auswahl verteilte Dateisysteme
- Ceph Theorie
- Ceph Praxis

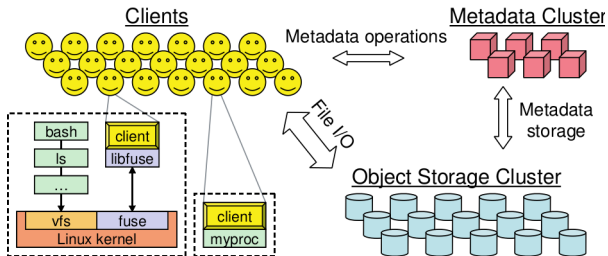
## 3 Erfahrungen

## Drei Komponenten

- **Client:** Posixähnliche Schnittstelle
- **OSD Cluster:** Speichert Daten und Metadaten
- **Metadaten Cluster**
  - Verwaltet Namensraum
  - Konsistenz/Kohärenz
  - Sicherheit (Noch nicht ausreichend implementiert)

## Drei Zugriffsmöglichkeiten

- Object-Store (Swift, Amazon S3, RESTful, ...)
- Block-device (RBD)
- Filesystem (CephFS, POSIX, Snapshots)



Quelle: [?]

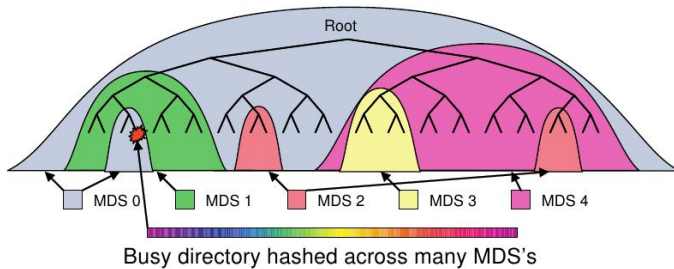
- Trennung der Verwaltung von Daten und Metadaten
  - Unabhängige Skalierbarkeit
  - Keine Block/Objektlisten (CRUSH)
- Dynamische Verteilung der Metadaten
  - Halber Aufwand sind Metadatenoperationen
  - Gutes Loadbalancing der Metadaten ist äußerst wichtig
- Autonomes OSD Cluster verantwortlich für:
  - Replikation
  - Fehlererkennung
  - Wiederherstellung nach Fehlern

# Dateizugriff (CephFS)

- ① Client sendet Dateirequest an das MDS Cluster
- ② Ein MDS übersetzt Dateiname in *file inode*
  - Unique inode number
  - Besitzer
  - Größe
  - ...
- ③ MDS gibt inode number, Größe und Verteilungsstrategie zurück
- ④ Direkter Zugriff von Client auf OSD



# Dynamische Metadaten-Verteilung

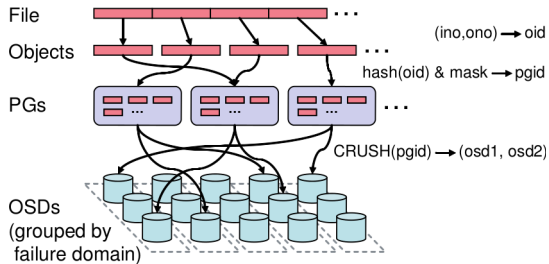


Quelle: [?]

**Problem:** Mehrere Clients schreiben/lesen eine Datei gleichzeitig.

- Kein read caching/write buffering mehr
- I/O der clients synchron
  - Synchrone I/O schlecht für Client Performanz
- HPC Erweiterung: O\_LAZY flag für open

# CRUSH



# CRUSH!

Für  $x \in \mathbb{N}$  ist *CRUSH* eine Abbildung nach  $R \in \mathbb{N}^n$  für  $n \in \mathbb{N}$ .

$x$  : Data-Object

$R$  : Placement für  $n$  Replicas

## 1 Einführung und Motivation

## 2 Parallel, distributed, fault-tolerant filesystems

- Auswahl verteilte Dateisysteme
- Ceph Theorie
- Ceph Praxis

## 3 Erfahrungen

- Ausgeschaltetes Write Caching auf der Festplatte
- Dateisystem mit *Extended Attributes (XATTRs)*:
  - XFS
  - BTRFS
- Empfohlen: Eine Platte je OSD

Installation einfach, Pakete vorhanden für:

- DEB (Ubuntu/Debian)
- RPM (CentOS, RHEL)

## 3 Daemons

- **ceph-osd**: Object Storage
- **ceph-mds**: Metadaten Server: Verteilter, kohärenter Cache für Metadaten
- **ceph-mon**: Monitor, Cluster Management, Konfiguration und Zugangspunkt für CephFS

# Object Storage Daemon

- Mehrere OSDs pro Host möglich
- Hardware Empfehlung:
  - Viele Festplatten
  - Viel Ram (besseres Caching) (min. 500MB/Daemon)
  - Schnelles Netzwerk
- Anzahl:
  - Mindestens 2, so viele wie möglich



- Hardware Empfehlung:
  - Sehr viel Ram
  - Schnelle CPU
  - Schnelles Netzwerk (geringe Verzögerung)
- Anzahl:
  - Mindestens 1
  - 2 oder mehr für Redundanz und Load Balancing

- Hardware Empfehlung:
  - Einige GB Festplattenspeicher
  - Feste IP
- Anzahl:
  - Ungerade Anzahl! (Split Brain)
  - 1 ist ok
  - 3 ideal für die meisten Anwendungsfälle
  - Mehr nur für sehr große Cluster

# Beispiel ceph.conf

```
1 [global]
2 auth supported = cephx
3 osd pool default size = 2 # Write an object 2 times.
4 [osd]
5 osd journal size = 1000
6 [mon.a]
7 host = myserver01
8 mon addr = 10.0.0.101:6789
9 [mon.b]
10 host = myserver02
11 mon addr = 10.0.0.102:6789
12 [mon.c]
13 host = myserver03
14 mon addr = 10.0.0.103:6789
15 [osd.0]
16 host = myserver01
17 [osd.1]
18 host = myserver02
19 [osd.2]
20 host = myserver03
21 [mds.a]
22 host = myserver01
```

- **ceph-deploy**

- ① Benötigt passwortlosen root-Login auf allen Hosts
- ② `ceph-deploy new mon`
- ③ `ceph-deploy osd create`
- ④ `ceph-deploy mds create`

# CRUSH-map

```
1 # buckets
2 host c096-137 {
3     id -6          # do not change unnecessarily
4     # weight 2.000
5     alg straw
6     hash 0 # rjenkins1
7     item osd.6 weight 1.000
8     item osd.7 weight 1.000
9 }
10 ...
11 rack unknownrack {
12     id -3          # do not change unnecessarily
13     # weight 6.000
14     alg straw
15     hash 0 # rjenkins1
16     item c096-134 weight 2.000
17     item c096-135 weight 2.000
18     item c096-136 weight 2.000
19     item c096-137 weight 2.000
20 }
21 root default {
22     id -1          # do not change unnecessarily
23     # weight 8.000
24     alg straw
25     hash 0 # rjenkins1
26     item unknownrack weight 8.000
27 }
```

## 1 Einführung und Motivation

## 2 Parallel, distributed, fault-tolerant filesystems

- Auswahl verteilte Dateisysteme
- Ceph Theorie
- Ceph Praxis

## 3 Erfahrungen


## Gut

- Selbstheilung sehr gut.
- Feature-Set sehr gut (insbesondere Zugriffsmöglichkeiten).

## Schlecht

- Crash bei mount mit dem Kernel Treiber auf OSD-Server (per Design).
- Performance bei vielen kleinen Dateien/Datenbanken noch relativ gering (Hardware? Laut Entwicklern ist auch CephFS noch nicht ganz so ausgereift).

# Alte Tests / alter Cluster

virt				
 Phoronix Test Suite	lxc-on-ext4	lxc-on-gluster	lxc-on-cephnbd-ext4	lxc-on-cephnbd-direct-ext4
Apache Benchmark	7747.08	1653.77	6721.80	7108.74
NGINX Benchmark	8322.34	584.71	7009.10	7818.15
PostgreSQL pgbench	227.18	41.52	6.60	11.31
PostMark	1383	114	1099	1154
IOzone	49.69	27.07	7.86	7.31
IOzone	1485.55	98.72	1141.47	1120.62
C-Ray	194.12	190.51	187.42	190.70
POV-Ray	1268	1280	1275	1304
PHPBench	31765	31692	31580	31791
GraphicsMagick	80	75	80	79
OpenSSL	41.30	41.55	41.63	41.60
Timed MAFFT Alignment	18.84	19.58	20.65	19.68
Himeno Benchmark	282.63	290.40	294.70	293.70
LAME MP3 Encoding	38.83	39.15	39.05	39.17
LZMA Compression	270.05	303.19	268.22	264.46
Stream	5468.09	4637.33	5268.77	5503.62
Stream	5373.22	4548.70	5132.91	5380.74
Stream	5791.56	4958.28	5415.82	5883.06
Stream	5867.09	5034.31	5481.42	5959.16
PHORONIX-TEST-SUITE.COM				

Benchmark: <https://rcs.num.math.uni-goettingen.de/virt-benchmark/composite.xml>



# Vorläufige neue Tests

**NFS** : 1 Node, 12 Platten

**NFS (MI)**: 1 Node, 8 Platten (neuer)

**CephFS**: 6 Nodes, 4 OSD-Nodes, 8 SAN-Platten, 2 Replica

Setup	read Kb/s	write Kb/s	read Kb/s (64)	write Kb/s (64)
NFS (MI)	78995	44613	8623	2162
NFS	316	164	790	223
CephFS	137613	18012	843	211

## Positiv

- Sehr vielversprechend und in aktiver Entwicklung.
- Team sehr offen und hilfsbereit.
- Fülle an Zugriffsmöglichkeiten und Skalierbarkeit.

## Negativ

- Einiges unvollständig implementiert (Sicherheit)
- Insbesondere CephFS noch nicht voll Produktionsreif

Fragen? Anregungen? Vorschläge?

