

Intern: Ceph

Kurzeinführung in die verteilte Storage-Lösung

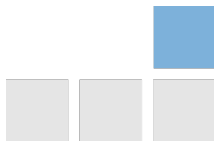


Dominik Vallendor ■ 29.05.2017



Motivation

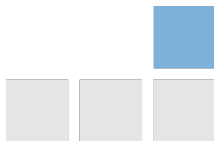
- Lokale Speicher sind unflexibel, selbst mit Redundanzlösungen (bsp. DRBD)
- Storages meist nicht redundant, teuer oder nicht erweiterbar
- iSCSI komplex zu handhaben
- Trend zu Cloud-Lösungen benötigt technisch soliden Unterbau



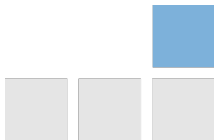


Was ist Ceph?

- Verteilte Storage-Lösung
- Unterstützt beliebige Anzahl von Servern
- Stark in die Breite skalierbar (Scale-Out)
- Objektspeicher mit S3-kompatibler Schnittstelle
- Virtuelles Blockdevice
- CephFS



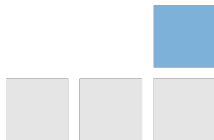
- Blockdevice sehr einfach einzubinden (rbd)
- Native Unterstützung des Blockdevices durch KVM/Libvirt
- Selbstheilend
- Kein Single Point of Failure
- Einfach erweiterbar mit neuen Festplatten und/oder Servern im laufenden Betrieb
- Redundanz-Anzahl frei wählbar und im Betrieb änderbar (meist 3)
- Upgrade-Pfade für Versions-Upgrade
- Updates von Nodes im laufenden Betrieb
- Snapshots und COW (ähnlich LVM) möglich



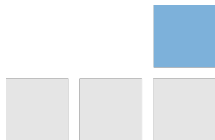
- Object Storage Nodes
- 1 bis x Object Storage Devices (OSD) pro Node
- Monitor Nodes
- Admin Node (optional)
- Separates Netz empfohlen

In der Regel:

mind. drei Server mit einem oder mehreren OSDs und jeweils einem Monitor



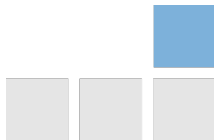
- *cephadmin*-Benutzer auf jedem Server
- Verwaltung vom Admin-Node aus, Zugriff per SSH (public key auth, Sudo)
- Verwaltung erfolgt meist durch *ceph-deploy*
- *ceph-deploy* ruft betriebssystem-spezifische Kommandos auf
- Software-Installation, starten von Diensten, etc. alles automatisch





Object Storage Device (OSD)

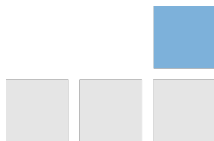
- liegen jeweils auf einer Partition
- *ceph-deploy* geht von direktem Festplatten-Zugriff aus
- LVM unerwünscht
- RAID unerwünscht und macht nur wenig Sinn
- XFS-Dateisystem auf der Partition
- Dateisystem wird im Storage Node lokal gemountet
- Verwaltung durch *ceph-osd* Daemon



Monitoring Node

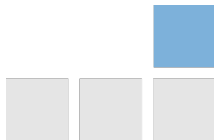
- Überwacht den Cluster
- Wird durch *ceph-mon* Daemon bereit gestellt
- Quorum erfordert mind. drei Server
- Läuft i.d.R. mit auf dem Storage Node (Doppelfunktion)

- RADOS (reliable autonomic distributed object store)
- Storage Pools
- Pool ID und Object ID
- Cluster-Map
- Hashing der Object ID und gleichmäßige Verteilung auf die Placement Groups
- CRUSH Algorithmus (Controlled Replication Under Scalable Hashing)



- Funktionsweise testen: *ceph health*
- Detaillierte Informationen: *ceph status*
- OSD-Status abfragen: *ceph osd tree*
- Umgang mit OSDs: *ceph osd out, stop, crush, del, rm, create*

- Messergebnisse noch nicht repräsentativ
- Testsystem mit langsamer Hardware, Partitionen schlecht verteilt
- Performance scheint insgesamt aber gut zu sein
- Zugriffsgeschwindigkeit ca. 30% bis über 100% einer einzelnen Festplatte
- Geschwindigkeit ist auch beim Hintergrund-Rebuild immer noch gut
- Detaillierte Übertragungs/Zugriffsdaten live einsehbar
- Performance angeblich stark abhängig vom Netzwerk. 10 Gbit/s empfohlen

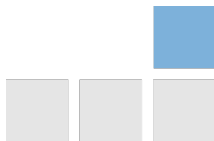


Storage Node Ausfall:

- Server wird aus dem Cluster entfernt
- Reaktion sehr schnell

Festplatten-Ausfall:

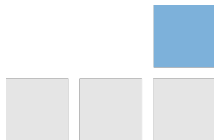
- Ausfall wird vom Storage Node erkannt
- OSD wird erst beim Zugriff entfernt
- Verzeichnis muss manuell geunmounted werden, etc.
- Prinzipiell aber kein Problem





OSD Replacement, Erweiterung

- Austausch eines defekten OSDs nicht vorgesehen
- Austausch = Entfernen + Erweitern von/mit OSDs
- Defekte OSDs müssen manuell deaktiviert/entfernt werden
- Cluster organisiert sich nach Änderung von OSDs automatisch neu



- Weitere Performance-Tests notwendig
- SSD-only vs. SSD-Journal
- Verschiedene Netzwerksetups und -bandbreiten testen
- Wie viele OSDs pro Storage Node?
- CRUSH-Map verstehen/bearbeiten
- KVM-Livemigration

