

2010 SourceTalk Tage

Hochperformante parallele Filesysteme für HPC

Andreas Landhäußer

T-Systems Solutions for Research GmbH

<mailto:Andreas.Landhaeusser@t-systems-sfr.com>

Mehr IOPS?

- mehr IOPS, als die einer einzelnen Festplatte, lassen sich nur sehr schwer erreichen.
- Lösung: mehrere Festplatten/Controller/IO-Karten werden in einem Verbund benutzt.
 - höhere Leistung als mit einer einzelnen Festplatte
 - mehrere Festplatten/Controller/IO-Karten
 - >> teurer, viel ungenutzter Plattenplatz.

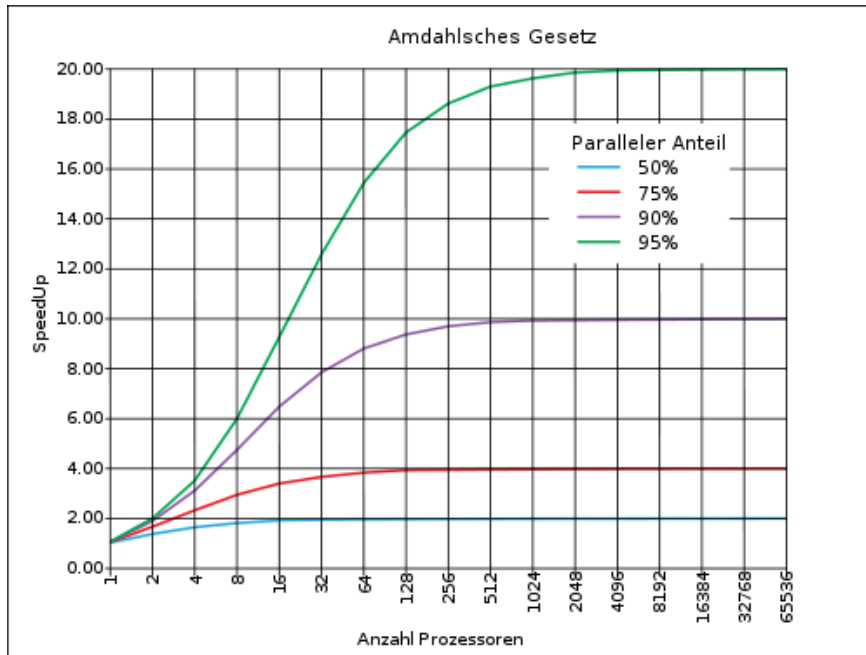
Redundant Array of Independent Disks (RAID)

- Erzeugung eines logischen Laufwerks aus mehreren Festplatten
 - RAID Controller
 - Software RAID
- Gebräuchliche RAID-Level
 - RAID 0 Striping = Beschleunigung ohne Redundanz
 - RAID 1 Mirroring = Spiegelung
 - RAID 5 Leistungssteigerung und Ausfallsicherheit

Amdahlsches Gesetz: Speedup

- $Speedup = \frac{1}{\sum_{k=0}^n \left(\frac{P_k}{S_k} \right)}$
- P_k Prozentanteil des Codes, der beschleunigt oder verlangsammt wurde $\sum P_k = 100$
- S_k Geschwindigkeitsmultiplikator (1 = keine Veränderung)
- k laufende Nummer an der eine Veränderung durchgeführt wurde
- n Gesamtanzahl der Änderungen

Amdahlsches Gesetz



$$S = \frac{1}{(1-P) + \frac{P}{N}} \leq \frac{1}{1-P}$$

Zugriffsmethoden auf Festplatten

- Network File System (nfs) Unix/Linux
- Common Internet File System (CIFS) Windows

Ermöglichen einen transparenten Zugriff von jedem Node auf das gemeinsame File System.

- Ohne Modifikation des Quellcodes nutzbar
- **Aber:** keine Abstimmung zwischen den Knoten, keine Cache-Coherence

Network attached Storage (NAS)

- Robust, da nur atomare Transaktionen

Nachteil:

- führt zu signifikanten Leistungseinbußen, da Protokoll Overhead mit steigender Netzleistung zunimmt
 - nfs Gbit Ethernet nur 60% max Leistung
 - CIFS unter 50%

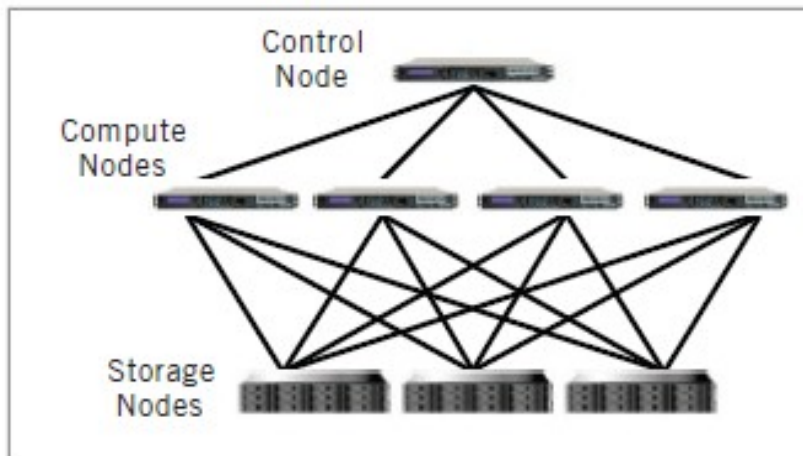
Storage Area Network (SAN)

- Leistungssteigerung und Vereinfachung des Managements über zentrale Speicher-Dienste
- Hochgeschwindigkeits Netzwerk zwischen SAN und Node
- Alle Knoten können auf das SAN zugreifen, aber ein gemeinsamer Zugriff ist nicht vorgesehen.

Parallele Filesysteme

Funktionweise:

Zentraler Kontrollknoten



Parallele Filesysteme

- Lustre
- GPFS
- pNFS
- FhGFS
- Ceph
- hadoop/HDFS
- ...

Lustre

- massiv paralleles Filesystem,
- 1998 an der CMU von Peter Braam vorgeschlagen,
- von CFS entwickelt,
- 2007 von Sun gekauft,
- 2009/2010 wurde Sun von Oracle erworben

Lustre: Kenndaten

- Single Name Space
- Skalierbarkeit: mehr Storage Server
- Performance: high Performance

Lustre: Bewertung

Vorteile:

- Skalierbar für 10000de Klienten, große Nutzerbasis
- High Performance
- Posix-Schnittstelle
- Stabilität

Nachteile:

- zentraler Metadatenserver, HA-Lösung möglich
- schwierige Administration
- wenig Zugriffssicherheit
- ungewisse Zukunft



Systems

GPFS: Kenngrößen

- Single Name Space
- Posix-Schnittstelle + non-Posix
- Skalierbarkeit: mehr Server, mehr JBODS
- Performance: high Performance
- Locking Schemata: File-, Block-, Byterange-level
- Zuverlässigkeit: Filereplication, Snapshots
- Administration: einfach während des Betriebs, running updates

GPFS: Bewertung

Vorteile

- Stabilität
- High Performance
- Skalierbar für 10000de Klienten, große Kundenbasis
- Posix-Schnittstelle
- verteilte Metadaten
- Einfache Administration

Nachteile

- kommerzielle Software
- Nur auf Linux, Windows und AIX portiert

pNFS

- NFS V1 – V4 sind nicht für Cluster geeignet
- pNFS ist eine Erweiterung von NFS V4 aka NFS V4.1
- Standardisierung erfolgte im Januar 2010

pNFS: Kenngrößen

- Separation von Kontroll- und Datenverkehrs
- NFS Server muss nicht mehr auch den Festplattenspeicher bereitstellen
- NFS Filesystem kann auf mehrere Speichersysteme verteilt werden
- Files können im Client, ähnlich wie bei den Block-Volume-Managern, repliziert oder auf mehrere Speichersysteme verteilt werden

pNFS: Bewertung

Vorteile

- nahtlose Integration ohne Änderungen
- abwärts kompatibel

Nachteile

- keine cache Kohärenz,
- Mehrere Metadatenserver, aber nur Kopien des Metadaten Bestands

FhGFS

- Entwickelt 2006 vom Institut ITWM an der FhG in Kaiserslautern
- Paralleles Filesystem für die HPC-Community
- FhGFS wird von FhG weiterentwickelt

FhGFS: Bewertung

- Vorteile
 - Freie Benutzung
 - Support durch Fraunhofer Gesellschaft
- Nachteile
 - kein HA
 - kein native Backup
 - zur Zeit nur Linux unterstützt

Ceph: Kenngrößen

- Nahtlose Skalierung, Ceph verteilt automatisch die Daten auf neue Storageeinheiten
- Alle Daten werden automatisch auf verschiedene Speichereinheiten repliziert, schnelle Wiederherstellung nach einem Fehler
- Adaptiver Metadatenserver, dynamische Anpassung an die aktuelle Workload
- Pseudo-random data distribution function (CRUSH)
- Reliable object storage service (RADOS)
- Extent B-tree object File System (today btrfs)

hadoop/HDFS: Bewertung

Vorteile:

- Basis für GRID Infrastruktur
- verteilte Daten (hoher Durchsatz)
- Archivierung
- heterogener Einsatz

Nachteile:

- kein general purpose Filesystem
- keine Posix-Schnittstelle
- wenig Zugriffssicherheit
- Java

