

Datentransfer zu HPC Clustern

Frühjarstagung AK Supercomputing 2024
Holger Angenent

Das Problem

- Port 22 von außerhalb der Uni in der Regel nicht (mehr) erreichbar
- Auf anderen Clustern könnte es schwierig werden, einen VPN-Client zu installieren
- Mit Jumphosts wird das vermutlich auch nicht schneller
- Man möchte die Daten nicht auf dem eigenen Rechner zwischenparken
- Nicht alle Use Cases abdeckbar
- Die Performance könnte suboptimal sein (obwohl man im Gegenzug auch nicht möchte, dass die Internetanbindung der Uni saturiert wird...)

Bisheriges Vorgehen

```
scp -r Verzeichnis Kennung@Cluster:
```

Mögliche Alternativen

- Globus Connect
 - Funktioniert (laut Leuten, die es verwenden)
 - GUI → für einige Leute einfacher
 - Closed Source
 - Teuer

Mögliche Alternativen

- Rucio
 - Entwickelt vom CERN
 - “Rucio enables centralized management of large volumes of data backed by many heterogeneous storage backends.”
 - Hat Kommandozeilenclient
 - Einigermaßen komplex

Mögliche Alternativen

- File Transfer System (FTS)
 - Entwickelt vom CERN
 - FTS3 is the service responsible for globally distributing the majority of the LHC data across the WLCG infrastructure. Is a low level data movement service, responsible for reliable bulk transfer of files from one site to another while allowing participating sites to control the network resource usage.
 - Kommandozeilenclient
 - Mir ist nicht genau klar, ob hier (einfach) einsetzbar

Ja hast du keine Tools auf Lager, die nicht vom CERN sind?

- “Wenn man als Werkzeug nur einen Hammer hat, sieht jedes Problem aus, wie ein Nagel” - Konfuzius



Mögliche Alternativen

- Nextcloud
 - Eigentlich kein Datentransfertools, sondern für Sync and Share
 - Open Source
 - Hat ein GUI
 - Zugang nicht über Port 22, sondern per 443
 - Auf Client Seite rclone (oder der Nextcloud Kommandozeilenclient) einsetzbar
 - Idealerweise hätte man eine rclone-Integration, so dass man zwei Nextclouds koppeln könnte (war mal in einem Projekt angedacht...)
 - Eröffnet weitere Use Cases wie das Empfangen von Daten von Externen → Bei Bedarf Demo:
<https://palma-web.uni-muenster.de>

rclone

- "The Swiss army knife of cloud storage"
- Rclone is a command-line program to manage files on cloud storage
- Transfers
 - MD5, SHA1 hashes are checked at all times for file integrity
 - Timestamps are preserved on files
 - Operations can be restarted at any time
 - Can be to and from network, e.g. two different cloud providers
 - Can use multi-threaded downloads to local disk
- Verwendung mit unserem Cluster siehe: <https://confluence.uni-muenster.de/display/HPC/Data+Transfers>

Messwerte Nextcloud

- Zwischen PC² Paderborn und Uni Münster
 - scp Paderborn -> Münster (10GB, home zu home): 59 MB/s
 - rsync Paderborn -> Münster (10GB, home zu home): 61 MB/s
 - rclone Paderborn -> Münster via Nextcloud (10GB, home zu home): 75 MB/s
 - scp Münster -> Paderborn (scratch zu home, 10 GB): 140 MB/s
 - rsync Münster -> Paderborn (scratch zu home, 10 GB): 138 MB/s
 - rclone Münster -> Paderborn via Nextcloud (scratch zu home, 10 GB): 420 MB/s

Nextcloud – technische Umsetzung

- Server mit Zugang zum Internet und zum parallelen Dateisystem (oder einem Knoten, der das kann) nötig
- Zugang von Nextcloud zum PFS per SMB oder SFTP möglich. (Hat beides Vor- und Nachteile)
- Braucht die Kennungen des Clusters
- Siehe: <https://palma-web.uni-muenster.de>
- Einbindung in Nextcloud per “External Storage” → Ansonsten verlangt Nextcloud exklusiven Speicherzugriff. So ist es möglich, dass Nextcloud auch die Daten erkennt, die nicht über Nextcloud geschrieben werden

Nextcloud – technische Umsetzung

- Voraussetzungen, Vor- und Nachteile beim Zugriff aufs Dateisystem per SMB
 - SMB und LDAP mögen sich nicht so, also muss der SMB-Server ins AD :(
 - Nextcloud muss sich mit Credentials gegen den Server authentifizieren
 - Da stehen ja nachher (gehashte) Nutzerpasswörter in der Datenbank!
- Dann macht man eben SFTP, da genügen public-/private Keys
 - Die Nutzenden müssen die einmal eintragen
 - Es läuft etwas langsamer

Nextcloud – weiteres Vorgehen

- “Dein Plan ist doof, weil die Nutzenden immer noch mit einem Kommandozeilentool hantieren müssen.”
- Das war auch ursprünglich anders geplant
- Eine Kopplung von zwei Nextclouds per rclone sollte das (u.a. vom CERN initiierte) Projekt CS3MESH4EOSC schaffen
- Vielleicht schaffen wir das auch noch, dass man über einen (Federated) Share Daten verschieben kann

Zusammenfassung

- Ich nextcloudifiziere mein paralleles Dateisystem
 - Dazu muss man ein paar Klümmzüge wegen des nicht-exklusiven Zugriffs/des IDMs machen
 - Port 443 nach außen exponieren (das machen wir bei anderen Projekten wie sciebo aber auch...)
- Auf der anderen Seite muss rclone verfügbar sein
- Würde jemand rclone irgendwie geschickt in Nextcloud integrieren, könnte man damit hoffentlich auch Daten direkt transferieren



Universität
Münster

**Vielen Dank für's Zuhören!
Fragen? Ideen?**