

# Austausch von Forschungsdaten zwischen Sites

Generische Tools in der HEP und Fragestellungen

*Oliver Freyermuth*

Universität Bonn, Physikalisches Institut  
freyermuth@physik.uni-bonn.de

3. Juli 2019

Datenerzeuger, Vor-Verarbeiter, Analyseteam und Archivierungsort sind selten „unter einem Hut“

**Beispiele:** CLS (*Coordinated Lattice Simulations*), HEP (*High Energy Physics*), Multi-Messenger-Astronomie, Machine-Learning. . .

## Fragestellungen

- Identitätsmanagement / Rechteverwaltung / Accounting
- Metadaten / Datensatz-Identifizier
- Automatisierte Replikation von Daten
- Einheitliche Schnittstellen
- Erkennung und Reparatur beschädigter Files
- Tape-Archivierung
- Analyse / Computing Power ohne lokales Storage
- Open Data

# Identitätsmanagement, Beispiel HEP

- Jeder Nutzer beantragt ein Zertifikat bei einer CA der IGTF (*Interoperable Global Trust Federation*)
- Nutzer registriert sich mit dem Zertifikat bei einer VO (*Virtual Organization*)
- Rollen / Capabilities / Untergruppen in einer VO möglich
- Regelmäßiges Unterzeichnen einer AUP (*Acceptable Use Policy*)
- Accountende: AUP nicht unterzeichnet / Zertifikat abgelaufen
- Anbindung IGTF an eduGAIN / DFN AAI

## Beispiel

```
DN: /C=DE/O=GermanGrid/OU=UniBonn/CN=Oliver Freyermuth
attribute: /atlas/de/Role=production/Capability=NULL
attribute: /atlas/de/Role=NULL/Capability=NULL
attribute: /atlas/Role=NULL/Capability=NULL
attribute: nickname = ofreyerm (atlas)
```

# Identitätsmanagement, Beispiel HEP

- Jeder Nutzer beantragt ein Zertifikat bei einer CA der IGTF (*Interoperable Global Trust Federation*)
- Nutzer registriert sich mit dem Zertifikat bei einer VO (*Virtual Organization*)
- Rollen / Capabilities / Untergruppen in einer VO möglich
- Regelmäßiges Unterzeichnen einer AUP (*Acceptable Use Policy*)
- Accountende: AUP nicht unterzeichnet / Zertifikat abgelaufen
- Anbindung IGTF an eduGAIN / DFN AAI

## Aber...

X.509 wird als Nutzerschnittstelle als Auslaufmodell gesehen!  
WLCG (Worldwide LHC Computing Grid) entwickelt sich zu  
Token-basierter Authentifizierung (SciTokens)

# Site-übergreifendes Handling von Datensätzen

- Datensätze haben innerhalb eines Experiments / VO einheitliches Namensschema (mit Scope)
- Dateien, Datensätze & Container, deren Speicherorte & Checksummen werden in einem zentralen Katalog erfasst
- spez. Zusatz-Kataloge (Analyseschritt, Version, Tags etc.)
- Nutzer reden mit Katalog-Server (*Rucio*)
  - Replikations-Regeln anlegen
  - Subscription über Namens-Pattern
- Rucio triggert Transfers über FTS-Server (*File Transfer Service*)
- Beschädigte Kopien werden erkannt, automatisch neu repliziert
- Rucio kennt angebotene Protokolle und Pfade zu den Dateien
- Tape-Archivierung und Staging über Speicher-Spacetoken (Ideen, dies transparent zu machen)

# Datenzugriff und Protokolle

- Klassisch: SRM / GridFTP (Spezial-Protokolle für Grid-Nutzung, außerhalb der HEP wenig Relevanz)
- Aktuell: Wandel zu HTTP (WebDAV) und XRootD, standardisiert, frei
- Protokolle müssen unterstützen:
  - Third-Party-Copy
  - Redirection / Ressourcen-Clustering
  - Delegation von Rechten, z.B. mit Token
  - Streaming von Daten (für Sites ohne lokalen Storage)
- Job-Routing zu den Daten möglich
- Transparentes Caching (XRootD)
- Idee von Data Lakes (regionale Storages)

## Fragestellungen

- Identitätsmanagement / Rechteverwaltung / Accounting
- Metadaten / Datensatz-Identifizier
- Automatisierte Replikation von Daten
- Einheitliche Schnittstellen
- Erkennung und Reparatur beschädigter Files
- Tape-Archivierung
- Analyse / Computing Power ohne lokales Storage
- Open Data

## Generische Tools, verbreitet in der HEP

- Rucio:  
<https://rucio.cern.ch/>
- File Transfer Service (FTS):  
<http://fts.web.cern.ch/>
- XRootD  
<http://xrootd.org/>



Danke

für die Aufmerksamkeit!

