

gLite Data Management Components

*Presenter : Ákos Frohner, CERN, IT-GD
On behalf of the Grid Data Management Team
EGEE'07, 1-5 Oct 2007, Budapest Hungary*

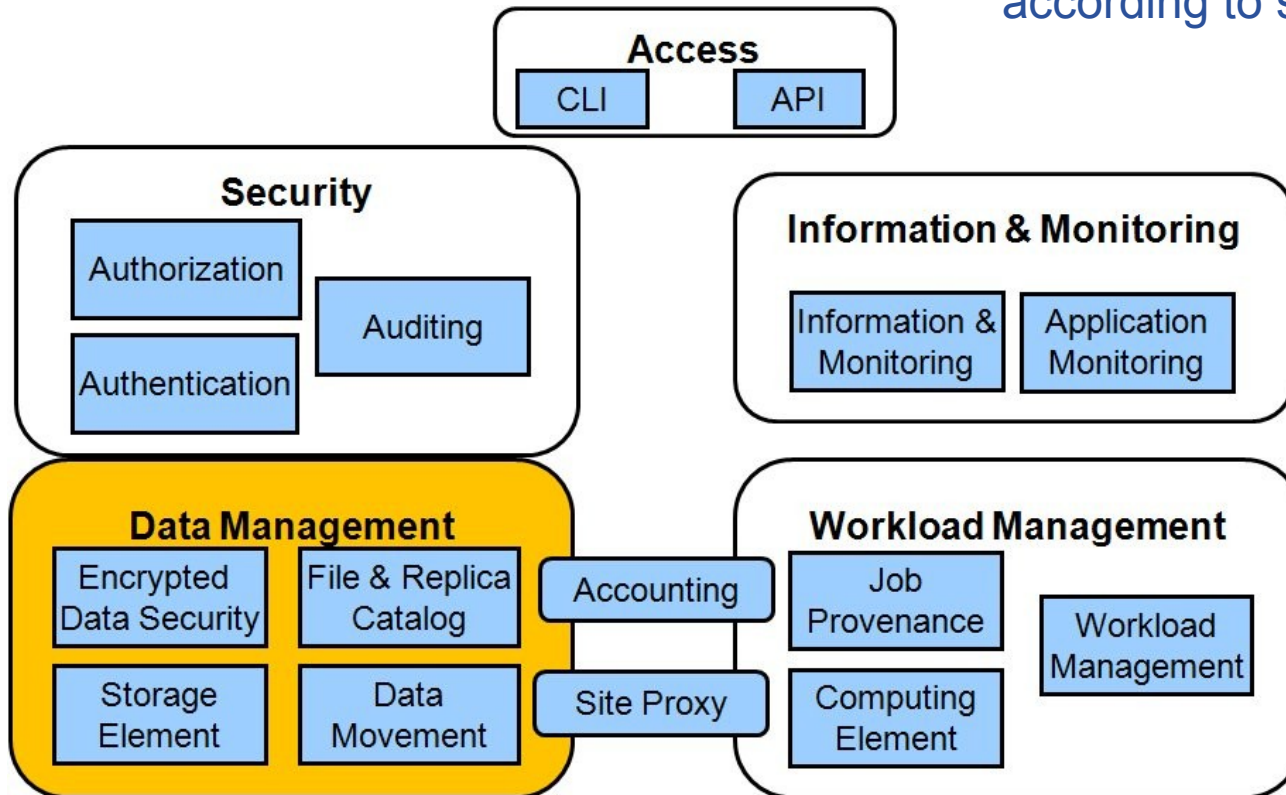


Akos.Frohner@cern.ch

- ***Will focus on EGEE data management contribution***
- **Enabling Grids for E-ScienceE (EGEE)**
 - Grid infrastructure for science (HEP, medicine, astronomy, ...)
 - 240+ sites, 45+ countries
 - Uses gLite as a lightweight, open source middleware distribution
- **Worldwide LHC Computing Grid (WLCG)**
 - Data processing based on a Tier-Model (Tier-0, Tier-1, Tier-2)
 - Use of Open Science Grid (OSG), EGEE (EGEE), NDGF, +
 - 15 PB/year to be stored at rates up to 1.5GB/sec (ALICE) and 100-150 MB/sec (ATLAS, CMS, LHCb)
 - Data sharing : ~500 Institutes, 5000 physicists, computer scientists and engineers

Service Oriented Architecture

- interoperability between grids
- support of grid standards
- flexible exploitation of the grid services according to specific needs



VO Frameworks

User Tools

**lcg_utils
FTS**

Data Management

GFAL

Cataloging

Storage

Data transfer

Information System/Environment Variables

Vendor Specific APIs

(RLS)

LFC

SRM

(Classic SE)

gridftp

RFIO

LFC

LCG File Catalog

LHC Computing Grid File Catalog

**Large Hadron Collider Computing Grid File
Catalog**

- **The LFC stores mappings between**
 - Users' file names and file locations on the Grid
 - Stores Permissions and
 - Ownership
 - Simple metadata

LFC file name 1

...

LFC file name n

"Replicas" are "Copies"

GUID

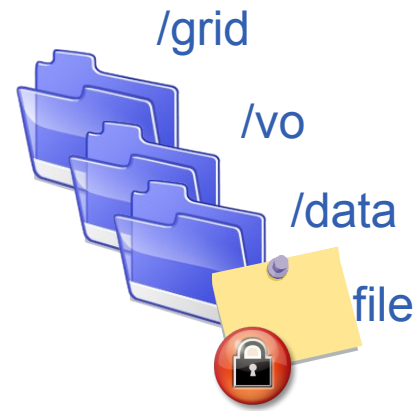
File replica 1

File replica 2

...

File replica m

- **Provides a hierarchical name space**
- **Supports GSI security model**
 - Including VOMS based ACLs
 - Very fine grained control
 - Implementation based on virtual IDs
 - Soon: encrypted channels
- **Simple DLI interface**
 - Data Location Interface
 - GUID <----> Location
 - Integration with WMS&RBs



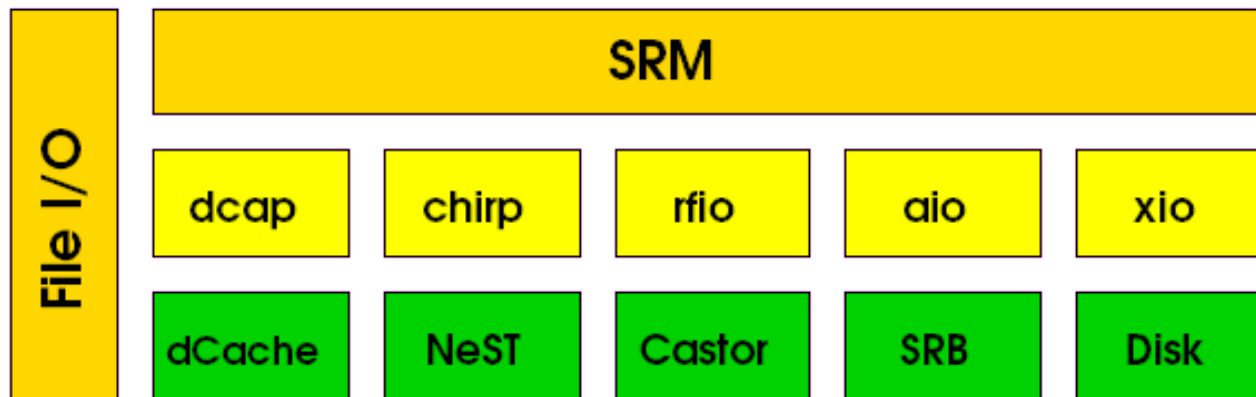
All files are "Write Once"

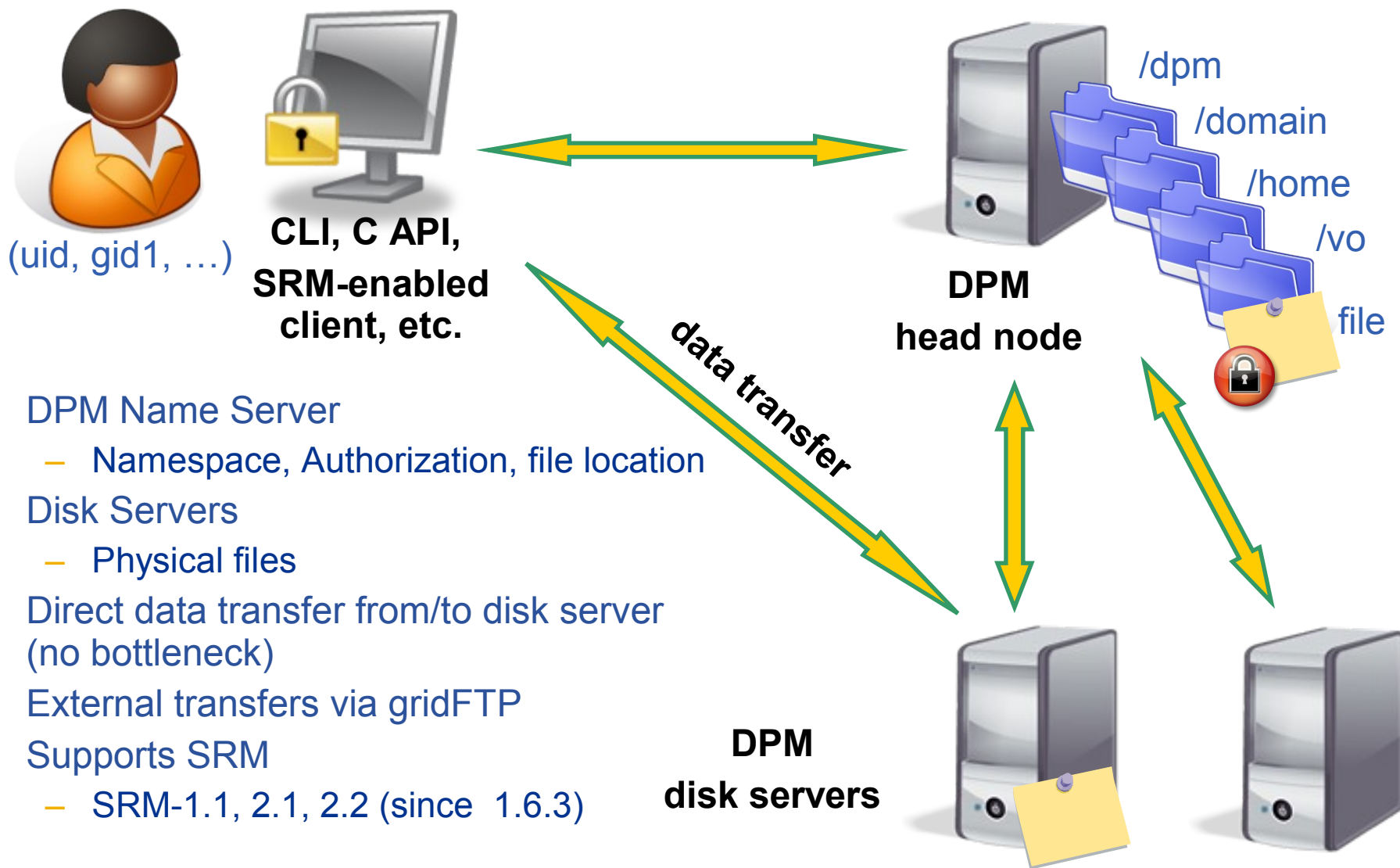
- **MySQL and ORACLE back-ends**
 - Ensures scalability and allows small scale deployment
 - Read only replication of catalogue (awaiting wider deployment)
- **Multi-threaded C server**
 - Supports multiple instances for load balancing
- **Thread-safe C clients**
 - Python & Perl bindings
 - Command line interface
- **Supports sessions to avoid authentication costs**
 - GSI is very expensive!
- **Bulk methods to reduce the number of round trips**
 - Under test by ATLAS --> 20 times faster
- **Widely used in EGEE:**
 - **largest LFC instance contains 8 millions entries**

DPM
Disk Pool Manager
SRM
Storage Resource Manager

- **Storage Resource Manager (SRM)**
 - Standard that hides the storage system implementation (disk or active tape)
 - handles authorization
 - Web service based on HTTPG
 - translates SURLs (Storage URL) to TURLs (Transfer URLs)
 - disk-based: DPM, dCache, Storm; tape-based: Castor, dCache
 - SRM-2.2
 - Space tokens (manage space by VO/USER), advanced authorization,
 - Better handling of directories, lifetime

- **File I/O: posix-like access from local nodes or the grid**
 - GFAL (Grid File Access Layer)

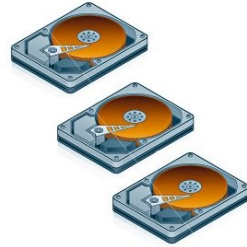




- DPM Name Server
 - Namespace, Authorization, file location
- Disk Servers
 - Physical files
- Direct data transfer from/to disk server (no bottleneck)
- External transfers via gridFTP
- Supports SRM
 - SRM-1.1, 2.1, 2.2 (since 1.6.3)

Addresses the storage needs of Tier-2 and smaller sites

- Focus on easy setup and maintenance
- Multi-threaded C implementation
- Name server DB
 - Keeps track of the status of files and their physical locations
 - MySQL and ORACLE back ends
 - Simplifies integration in existing local DB infrastructure
 - Ensures scalability
 - Shares code with LFC --> fix once run twice!
- Thread-safe C client and command line interface
 - http/https DPM browser (implemented, very soon to be released)
 - users and site managers interact with DPM at different levels
- GSI and VOMS based authorization and fine grained ACLs
 - Implemented via virtual IDs -> no excessive use of pool accounts
 - Pool access control on VO basis



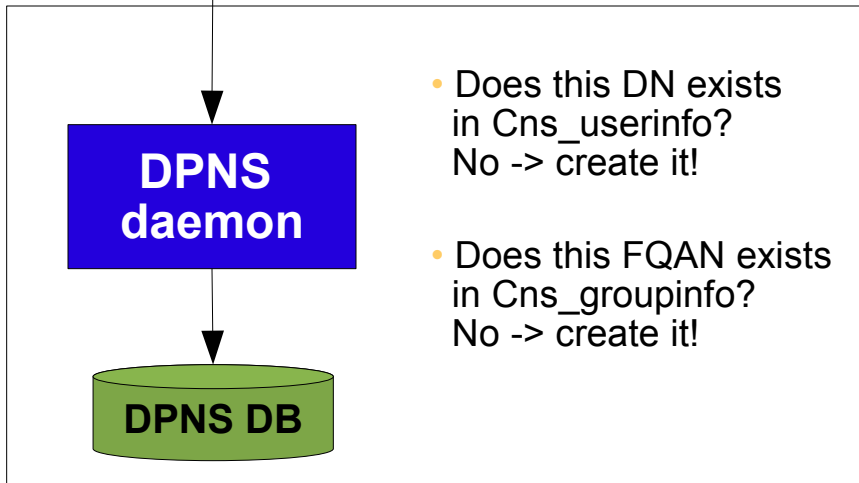
Different file access and transfer protocols

- Secure Remote File Input/Output (RFIO)
 - Secure file transfer and manipulation.
 - Implementation of thread-safe C client and a command line interface
 - Support of streaming mode.
- GSIFTP allows remote file transfer
 - New gridftp plugin is implemented to support gridFTP-2
- Xrootd: usable but still limited
 - no support of grid/voms certificates yet
- https/http: web access based on Apache.
 - Protocol http or https can be specified at transfer time

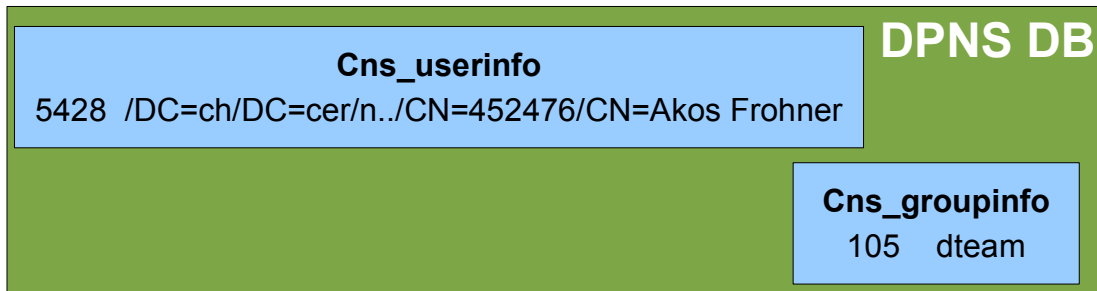
DN: /DC=ch/DC=cern/OU=Organic Units/OU=Users/CN=szamsu/CN=452476/CN=Akos Frohner

```
voms-proxy-init -voms dteam
```

```
dpns-ls /dpm/cern.ch/dteam/generated
```



- no need to create pool accounts
- no need to change the /etc/passwd file
- faster check on ACL than with string/pattern matching on DN/FQAN



- **EGEE Catalog**

- 110 LFCs in production
 - 37 central LFCs
 - 73 local LFCs

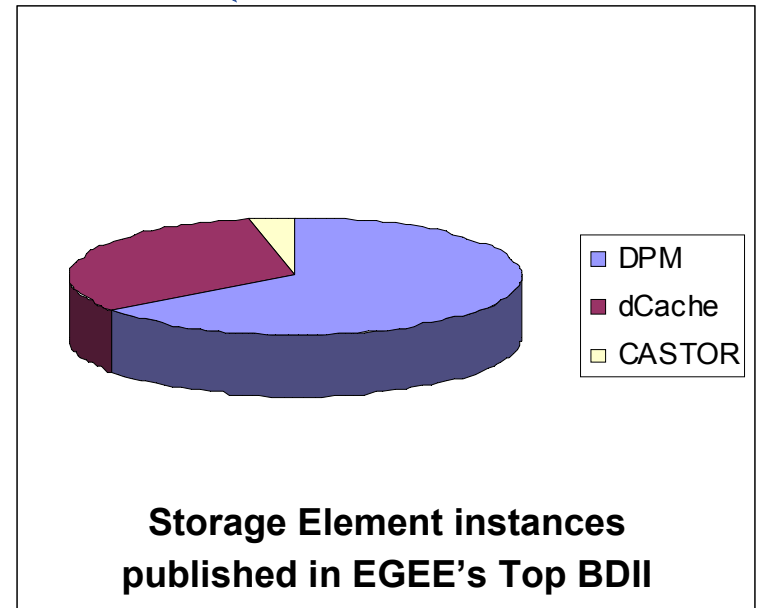
- **EGEE Storage Elements**

- CASTOR
- dCache
- DPM
 - 96 DPMs in production
 - Supporting 135 VOs

- **LFC and DPM**

- Stable and reliable production quality services
- Well established services
- Require low support effort from administrators and developers

Data volume distribution looks quite different



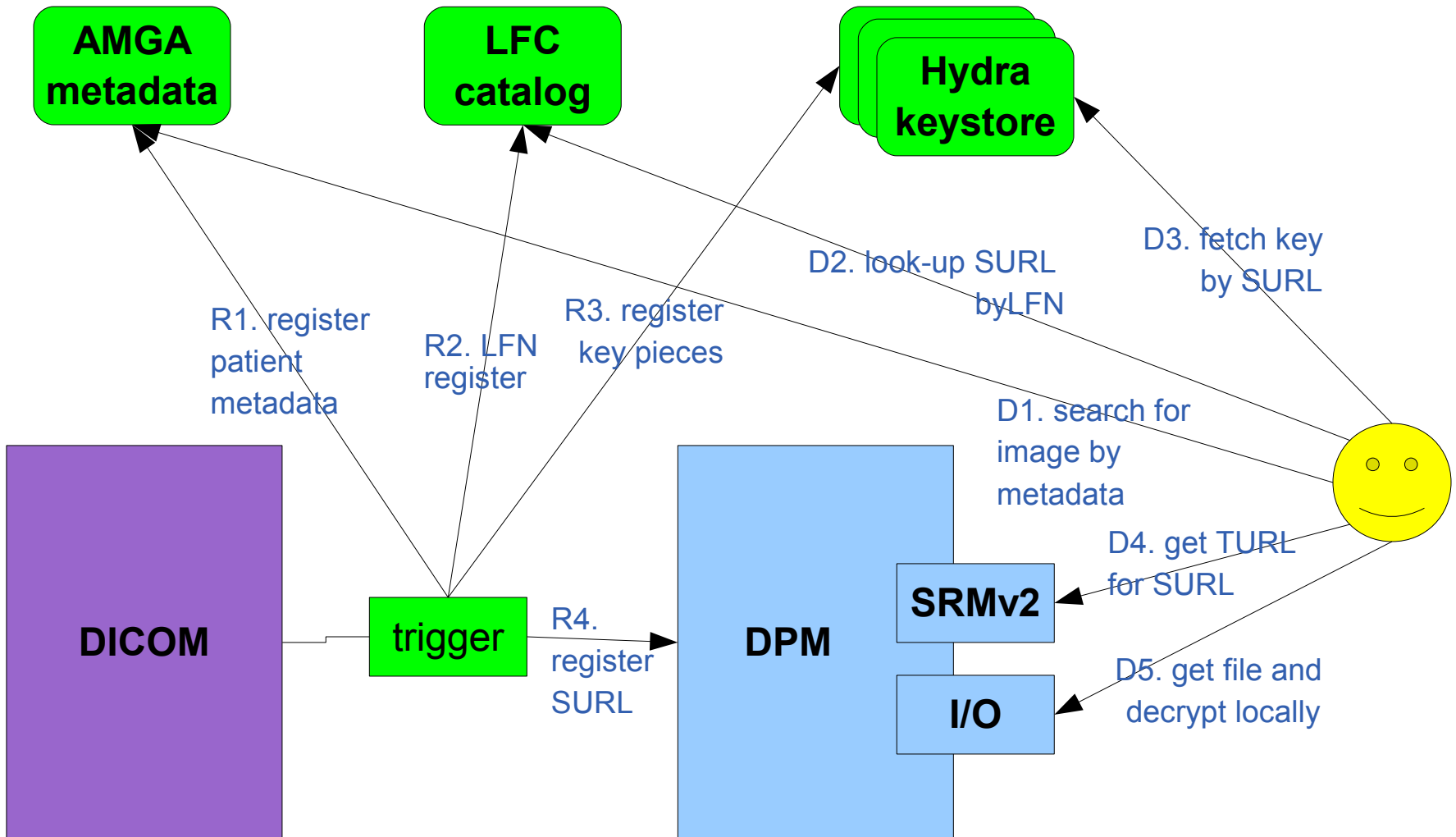
Problem : Medical institutes request data storage encryption

- Use of the DICOM standard for medical image handling
- Image retrieval and storage from/in DICOM servers : security issues

Solution : Extension of the data management tools (under way)

- File encryption on the fly, local decryption
- Use of HYDRA for split key management
- Use of the LFC to register/retrieve system data
 - Replicas location, filesize, ...
- Use of srmv2 to get the turls
- Use of I/O protocols, gridftp to load medical images
- Access control based on VOMS

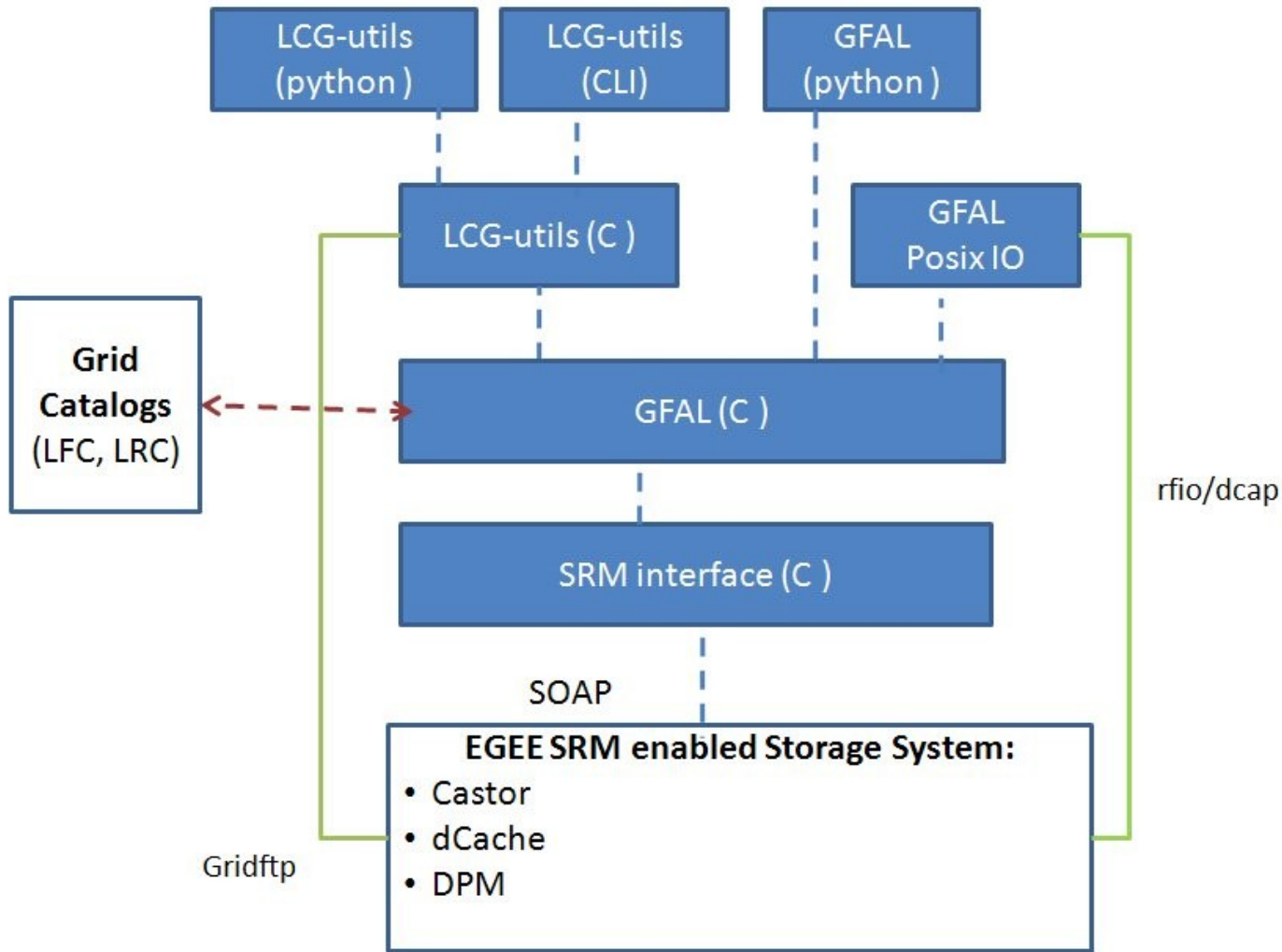




Lcg-util & GFAL

Grid File Access Layer

- **Purpose: Create the illusion of POSIX I/O**
 - Shield users from complexity
 - Interact with the information system, catalogue, SRMs
 - Can be used with/without information system/ catalogue
- **LCG-util :**
 - Command line and C-API
 - Covers most common use cases
 - Replication, catalogue interaction etc,
 - high level tool box
- **Gfal:**
 - Posix like C API for file access
 - SRMv2.2 support
 - user space tokens for retention policy (custodial/replica) & access latency (online/nearline)

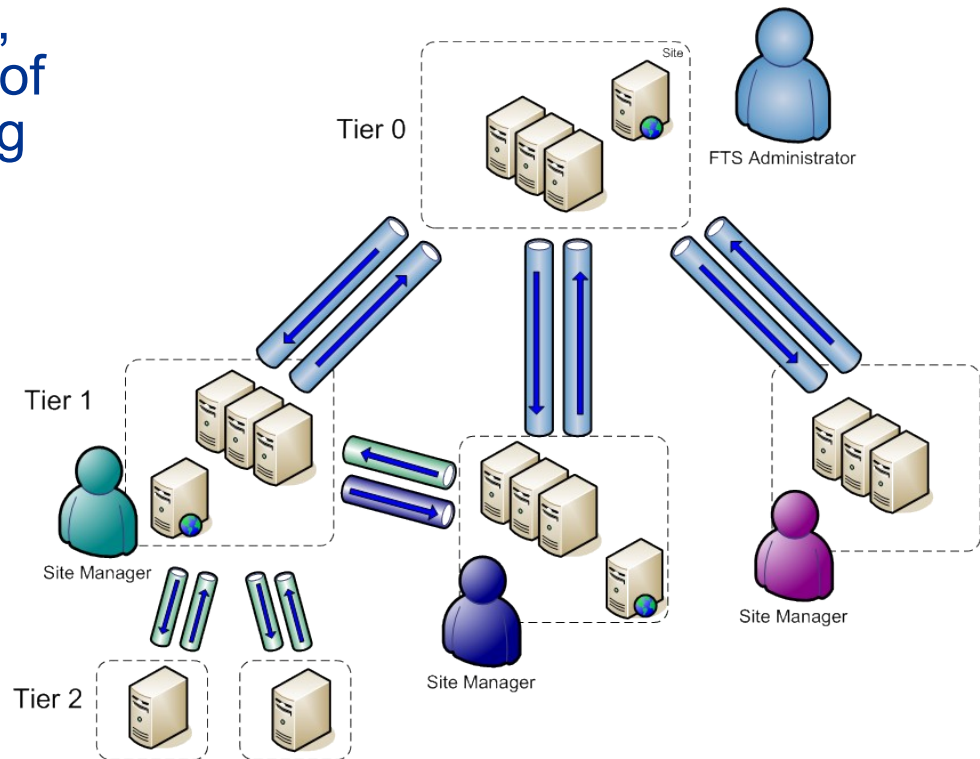


— File transfer protocol compatible with all EGEE storage systems

FTS File Transfer Service



- **gLite File Transfer Service is a reliable data movement service (batch for file transfers)**
 - FTS performs bulk file transfers between multiple sites
 - Transfers are made between any SRM-compliant storage elements (both SRM 1.1 and 2.2 supported)
 - It is a **multi-VO** service, used to balance usage of site resources according to the SLAs agreed between a site and the VOs it supports
 - VOMS aware



- **Why is it needed ?**

- For the **user**, the service it provides is the reliable point to point movement of Storage URLs (SURLs) and ensures you get your share of the sites' resources
- For the **site manager**, it provides a reliable and manageable way of serving file movement requests from their VOs and an easy way to discover problems with the overall service delivered to the users
- For the **VO production manager**, it provides ability to control requests coming from his users
 - Re-ordering, prioritization,...
- **The focus is on the “service” delivered to the user**
 - It makes it easy to do these things well with minimal manpower

- **Reliability**

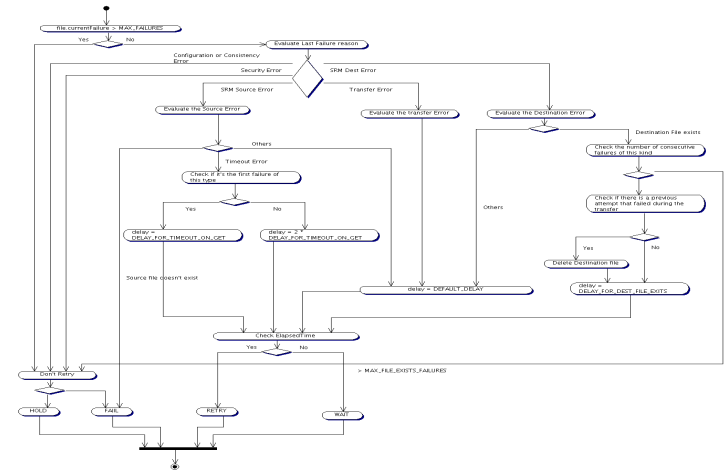
- It handles the retries in case of storage / network failures
 - VO customizable retry logic
- Service designed for high-availability deployment

- **Security**

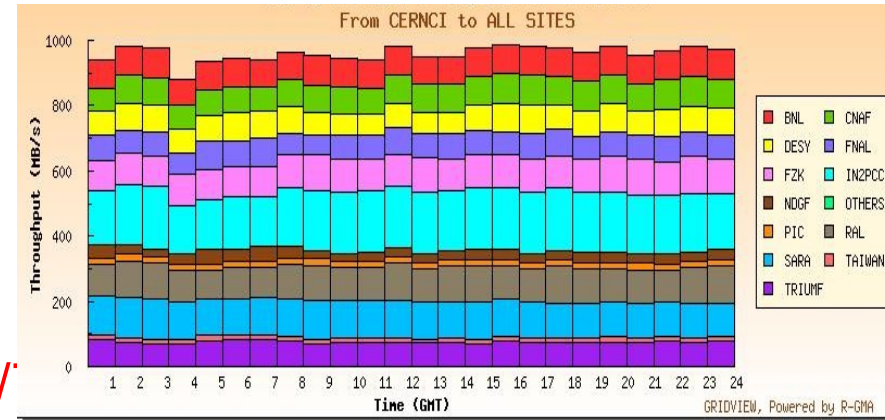
- All data is transferred securely using delegated credentials with SRM / gridFTP
- Service audits all user / admin operations

- **Service and performance**

- Service stability: it is designed to efficiently use the available storage and network resources without overloading them
- Service recovery: integration of monitoring to detect service-level degradation



- Designed to scale up to the transfer needs of very data intensive applications
- Currently deployed in production at CERN
 - Running the production WLCG tier-0 data export
 - Target rate is **~1 Gbyte/sec 24/7**
 - Over **9 petabytes** transferred in last 6 months > **10 million files**
- Also deployed at **~10 tier-1 sites** running a mesh of transfers across WLCG
 - Inter-tier1 and tier-1 to tier-2 transfers
 - Each tier-1 has transferred around 0.2 – 0.5 petabytes of data



- **New features in FTS 2.0**

- Better security model (certificate delegation)
- Support for SRM v2
- More administrative tools, more advanced monitoring features to make it easier to operate the overall service
- Soon:
 - Better support for clouds and channel sets
 - Black/White-listing SEs
 - Better integration with VO workflow management
 - *Call backs, hooks in the state machine*
 -

- **Focus continues upon service monitoring and easing the service operations together with closer integration of FTS with experiment software frameworks**

Definitely NO!!!

- The AMGA meta data catalogue (by Birger Koblitz)
 - Widely used by experiments
 - Is in the process to be integrated in the gLite distribution
- Many data management tools and services developed by Vos
- Lessons learned
 - A DM stack can only be developed with production feedback
 - The right balance between exposing details and hiding is hard to find
 - There will be more to do

Current status

- Data Management framework is usable
- LFC, FTS, DPM and lcg-util/gfal are used in production on a large scale

Outlook & Future

- ACL synchronization between LFC and SEs
- Improvements to lcg-util/gfal
 - e.g. flexibility to work independently of the LFC
- Better tools to check consistency in DPM
- Extension of Xrootd to support grid/voms certificates
- Finish medical data management implementation
- DPM : quota on pools and accounting
- Operational improvements to the FTS

- Continue the dialog with the user communities to focus effort