



Enabling Grids for E-science

# Data Management Services

*Riccardo Bruno*

*INFN*

*gLite Tutorial at the First EGEE User Forum*

*CERN, 27-28.02.2006*

[www.eu-egee.org](http://www.eu-egee.org)



Information Society



- **Grid Data Management Challenge**
- **Storage Elements, SRM**
- **LFC (LCG File Catalog)**
- **File and Replica Catalog (FiReMan)**
- **gLite I/O**

- **Grid Data Management Challenge**
- **Storage Elements (SRM)**
- **LFC (LCG File Catalog)**
- **File and Replica Catalog (FiReMan)**
- **gLite I/O**

# The Grid DM Challenge

NEEDS	REQUIREMENTS	SOLUTIONS
<p><b>Heterogeneous</b> : Data are stored on different storage systems using different technologies.</p>	<p>A common interface to storage resources is required in order to hide the underlying complexity.</p>	<p><b>Storage Resource Manager (SRM)</b> interface; (gLite File I/O Server)</p>
<p><b>Distributed</b>: Data are stored in different locations; in most cases there is no shared file system or common namespace.</p>	<p>Data need to be moved between different locations.</p> <p>There is need to keep track where data is stored.</p>	<p><b>File Transfer Service (FTS)</b> – to move files among GRID sites.</p> <p><b>Catalog</b> – to keep track where data are stored.</p>
<p><b>Data Retrieving</b>: Applications are located in different places from where data are stored.</p>	<p>There is need of scheduled reliable file transfer service.</p>	<p><b>File Transfer Service</b></p> <ul style="list-style-type: none"> <li>•Data Scheduler</li> <li>•File Placement Service</li> <li>•Transfer Agent</li> <li>•File Transfer Library</li> </ul>
<p><b>Security</b>: Data must be managed according to the VO membership access control policy.</p>	<p>Centralized Access control Service.</p>	<p><b>File Authorization Service</b></p>

- **DM** works with **files**, this assumption is due the following reasons:
  - **semantic** of file is very good understood by everyone
  - file is the **smallest granularity** of data.
- **EGEE's Specific Grid Requirements:**
  - **HEP** - High Energy Physics
  - **Biomed**

- **File Access Patterns:**
  - Write once, read many
  - Rare append - only updates with one owner
  - Frequently updated at one source - replicas check/pull new version
  - (NOT frequent updates, many users, many sites)
  
- **File naming**
  - Mostly, see the “logical file name” (LFN)
  - LFN must be unique:
    - includes logical directory name
    - in a VO namespace
  - E.g. /gLite/myVOname.org/runs/12aug05/data1.res

- **Storage Element – common interface to storage**
  - **S**Storage **R**esource **M**anager      Castor, dCache, DPM, ...
  - POSIX-I/O      gLite-I/O
  - Native Access protocols      rfio, dcap
  - Transfer protocols      gsiftp
  
- **Catalogs – keep track where data are stored**
  - File Catalog
  - Replica Catalog
  - File Authorization Service
  - Metadata Catalog

gLite File and Replica Catalog

**FireMan**

**LFC**

**AMGA Metadata Catalogue**
  
- **File Transfer – schedules reliable file transfer**
  - Data Scheduler
  - File Transfer Service  
(manages physical transfers)
  - File Placement Service  
(FTS and catalog interaction in a transactional way)

*(only designs exist so far)*

**gLite FTS – lcg-rep**

- **Grid Data Management Challenge**
- **Storage Elements (SRM)**
- **LFC (LCG File Catalog)**
- **File and Replica Catalog (FiReMan)**
- **gLite I/O**



He is running a job which needs:

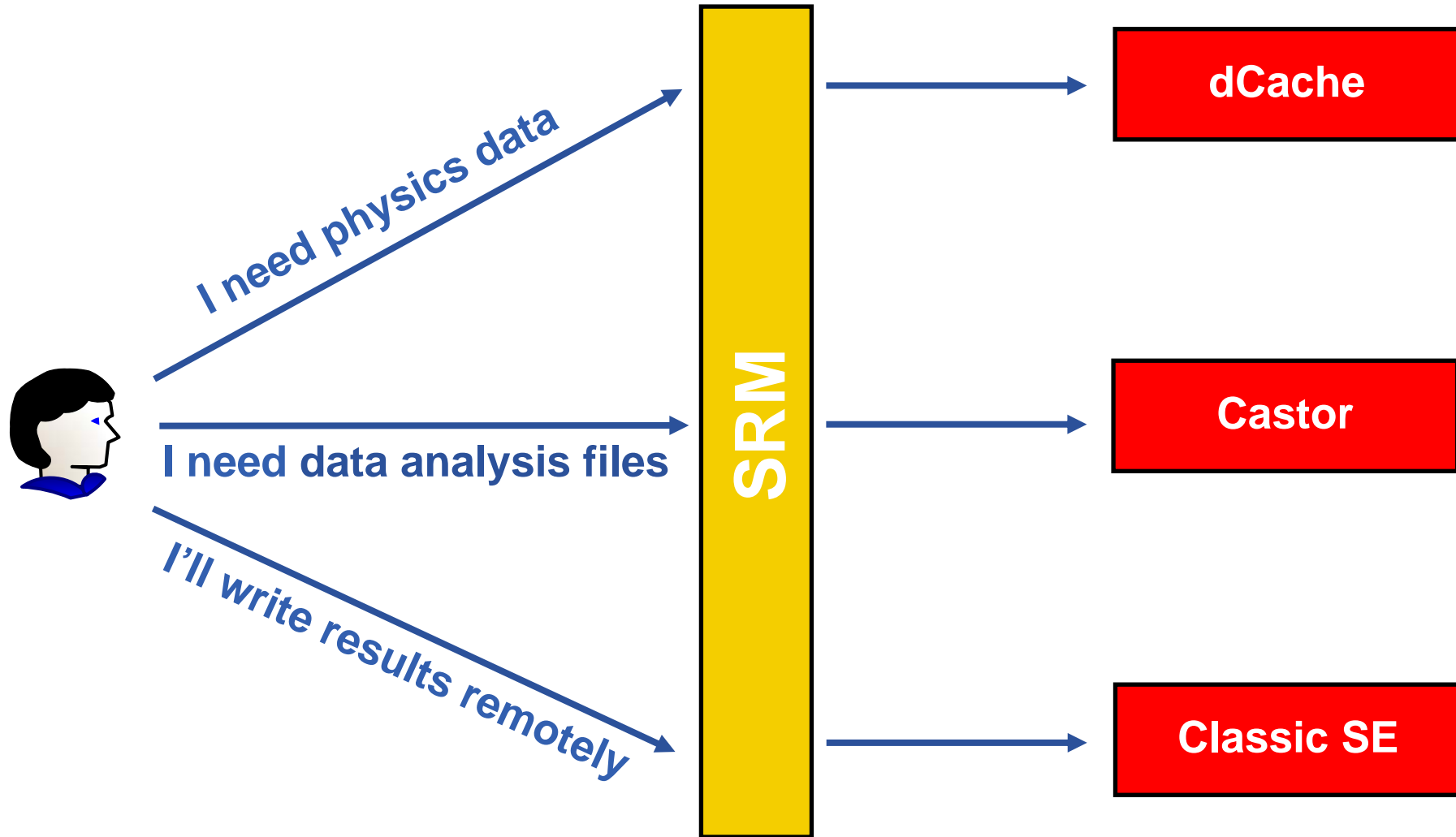
- Data for physics event reconstruction
- Some data analysis files
- He will write files remotely too

They are at CERN  
in dCache

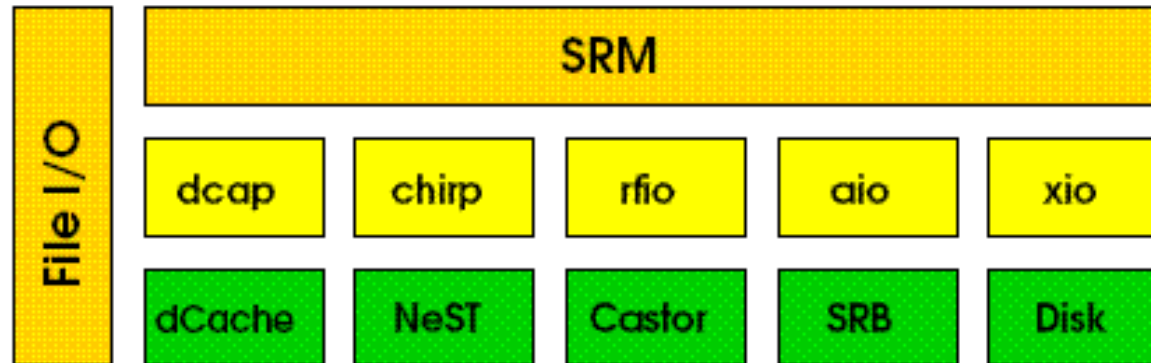
They are at Fermilab  
in a disk array

They are at Nikhef  
in a classic SE





- Data are stored on **Disk Pool Servers** or **Mass Storage Systems**
- **Storage Resource Management** needs to take into account
  - Transparent access to files (migration to/from disk pool)
  - Space reservation
  - File status notification
  - Life time management
- **SRM (Storage Resource Manager)** takes care of all these details
  - **SRM** is a **Grid Service** that takes care of **local storage** interaction and provides a **Grid interface** to outside world.
- Interaction with the **SRM** is hidden by higher level services

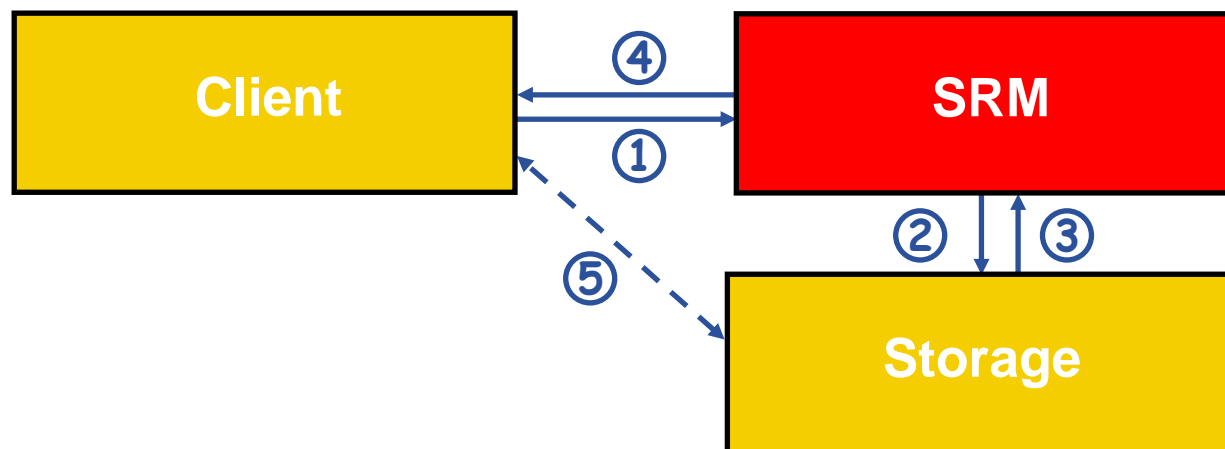


## Main supported MSS

- dCache
- NeST
- CASTOR
- SRB
- Disk

## Main supported protocols

- dCap
- Chirp
- Rfio
- Aio
- xio

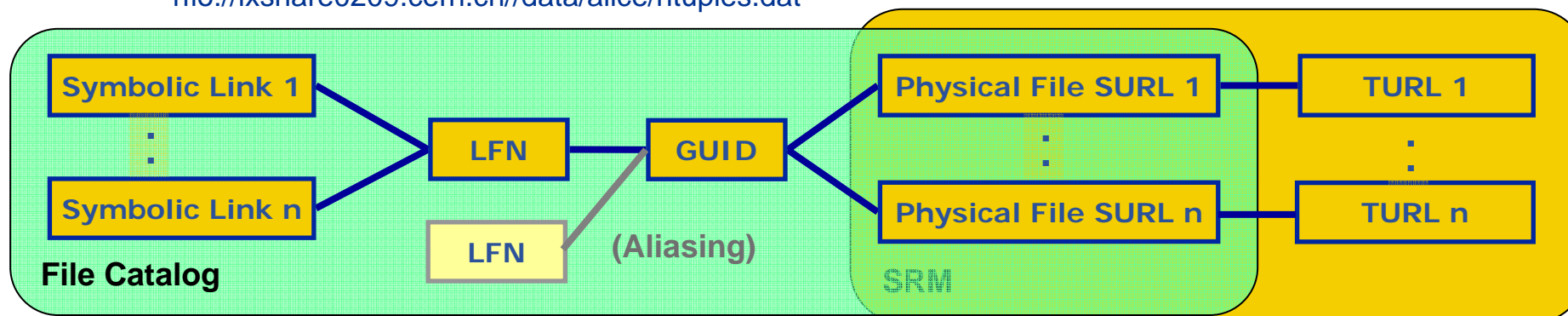


1. The client asks the SRM for the file providing a SURL (Site URL)
2. The SRM asks the storage system to provide the file
3. The storage system notifies the availability of the file and its location
4. The SRM returns a TURL (Transfer URL), i.e. the location from where the file can be accessed
5. The client interacts with the storage using the protocol specified in the TURL

- **Provide an SRM interface**
  - Specific Storage Solution: HPSS, CASTOR, DiskeXtender (UNITREE), DPM, dCache
- **Support basic file transfer protocols**
  - **GridFTP** mandatory
  - Others if available: https, ftp, etc.
- **Support a native I/O access protocol**
  - **POSIX** (like) I/O client library for direct access of data (rfio, dcap, gsidcap)
- **Other auxiliary services**
  - Accounting

- **Grid Data Management Challenge**
- **Storage Elements (SRM)**
- **LFC (LCG File Catalog)**
- **File and Replica Catalog (FiReMan)**
- **gLite I/O**

- **Symbolic Link** in logical filename space
- **Logical File Name (LFN)**
  - An alias created by a user to refer to some item of data, e.g. “lfn:cms/20030203/run2/track1”
- **Globally Unique Identifier (GUID)**
  - A non-human-readable unique identifier for an item of data, e.g. “guid:f81d4fae-7dec-11d0-a765-00a0c91e6bf6”
- **Site URL (SURL) (or Physical File Name (PFN) or Site FN)**
  - The location of an actual piece of data on a storage system, e.g.
    - “srm://pcrd24.cern.ch/flatfiles/cms/output10\_1” (SRM)
    - “sfn://lxshare0209.cern.ch/data/alice/ntuples.dat” (Classic SE)
- **Transport URL (TURL)**
  - Temporary locator of a replica + access protocol: understood by a SE, e.g. “rfio://lxshare0209.cern.ch//data/alice/ntuples.dat”





The **L**CG **F**ile **C**atalog fixes the performance and scalability problems of EDG (European Data Grid) file catalogs.

## Provides

- Bulk operations.
- Cursors for large queries.
- Timeouts and retries for client operations.

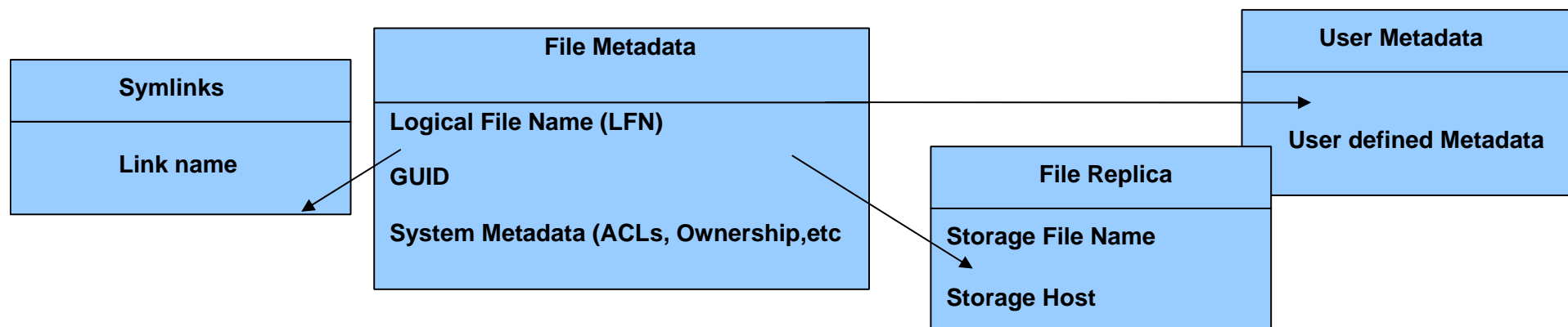
## Added features :

- User exposed transaction API.
- Hierarchical namespace and namespace operations.
- Integrated GSI Authentication and Authorization.
- Access Control Lists (Unix Permissions and POSIX ACLs).
- Checksums.

Supported database backends: **Oracle** and **MySQL**

GFAL integration and support to lcg-\* done by Grid Deployment group

- LFC stores both **logical** and **physical** mappings for the file in the same database → Speed up of operations
- Treats all entities as files in a **UNIX-like** filesystem.
- File **API** also similar to UNIX (create(), mkdir(), chown()....)
- Hierarchical namespace of **LFNs** mapped to the **GUIDs**
- **GUIDs** mapped to the physical locations of file **replicas** in the storage
- **System attributes** of files (creation time, file size and checksum...)  
stored as LFN attributes
- One field for **user-defined metadata**
- Multiple LFNs per GUID allowed as **symbolic links** to the primary LFN.

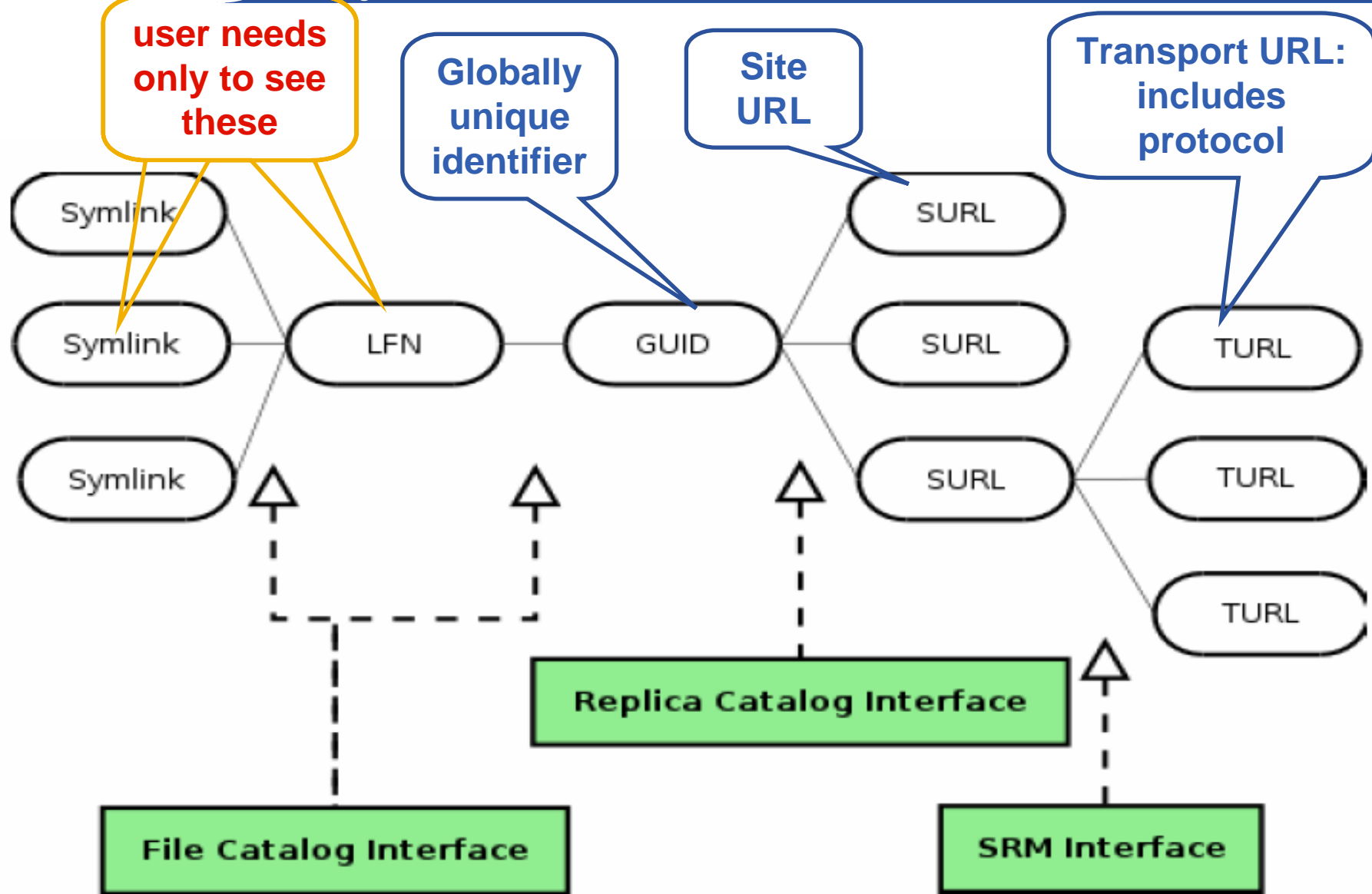


**If a site acts as a central catalog for several VOs, it can either have:**

- **One LFC server, with one DB account containing the entries of all the supported VOs. You should then create one directory per VO.**
- **Several LFC servers, having each a DB account containing the entries for a given VO.**

**Both scenarios have consequences on database backup policies.**

- **Grid Data Management Challenge**
- **Storage Elements (SRM)**
- **LFC (LCG File Catalog)**
- **File and Replica Catalog (FiReMan)**
- **gLite I/O**

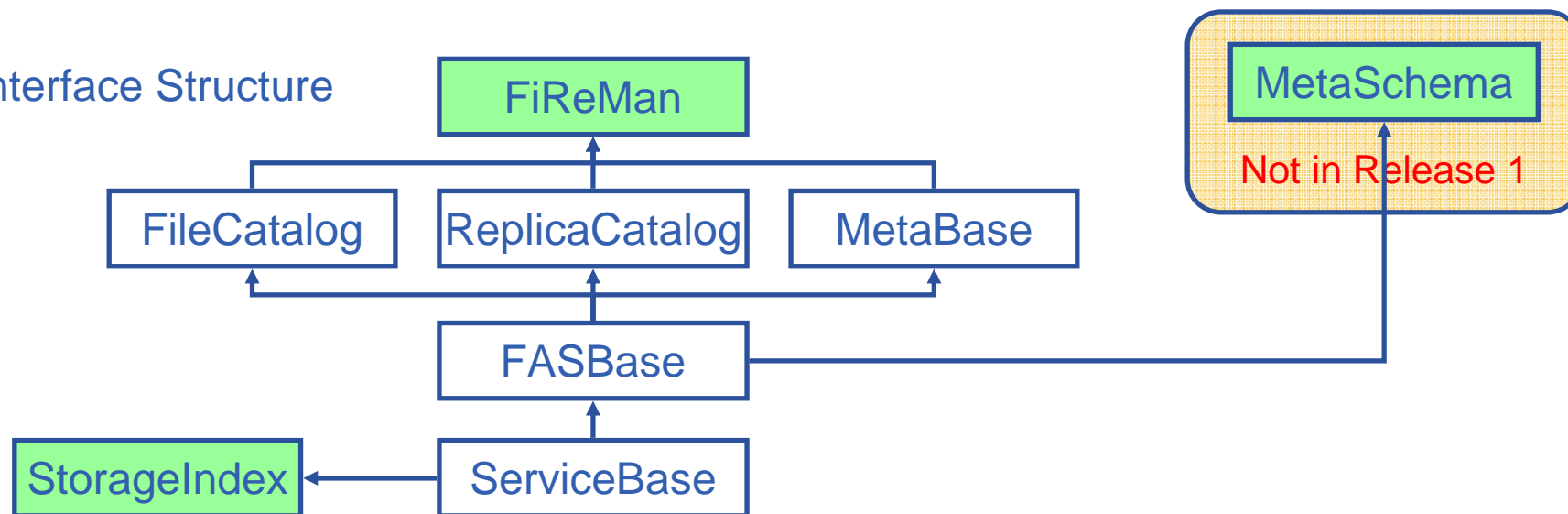


- **File Catalog**
  - Allows operations on the logical file namespaces that it manages (e.g.: making directories, renaming files, creating symbolic links)
  - Manages LFNs, keeping internally **LFN-GUID** mappings
- **Replica Catalog**
  - Exposes operations concerning the replication aspect of the grid files (e.g.: listing, adding and removing replicas of a file identified by its GUID)
  - Gives access to the **GUID-SURL** mappings
- **File Authorization Service (FAS)**
  - Request authorization - based on the DN and the Groups from the user's delegated credentials
- **StorageIndex**
  - Allows WMS interactions (file location for the RB)
- **Fireman = File and Replica Manager**
  - Provides all the previous services

- Logical File Namespace management
- Replica locations
- File-based metadata
- Metadata Management
- Authentication and Authorization information (ACLs)
- Service Metadata
- WMS interaction and global file location

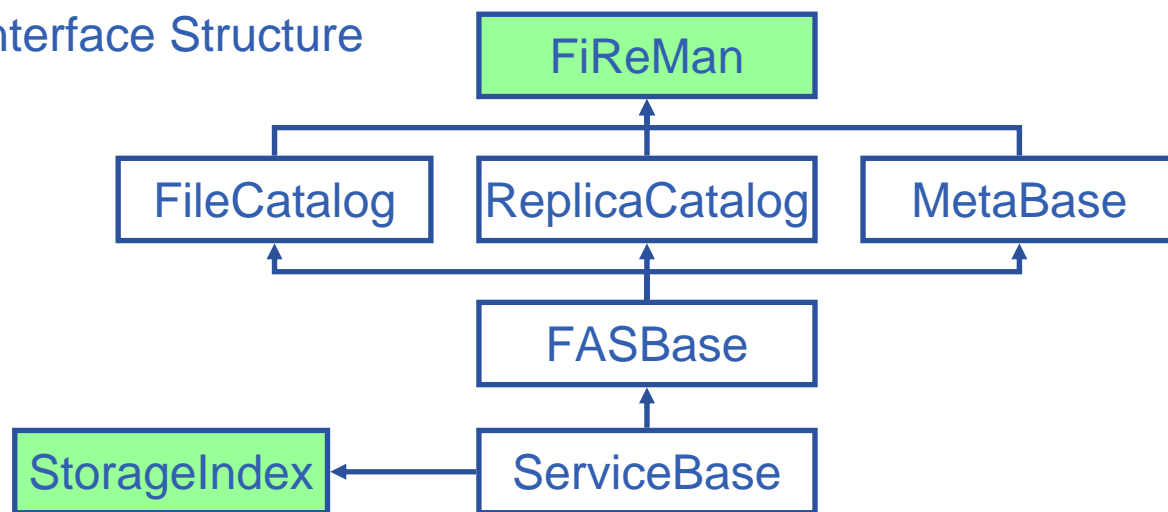
- FileCatalog**
- ReplicaCatalog**
- MetaBase**
- MetaSchema**
- FASBase**
- ServiceBase**
- StorageIndex**

Interface Structure



- Web Service interface (WSDL)
- Bulk operations
- Stateless interaction
- No transactions outside Bulk

## Interface Structure



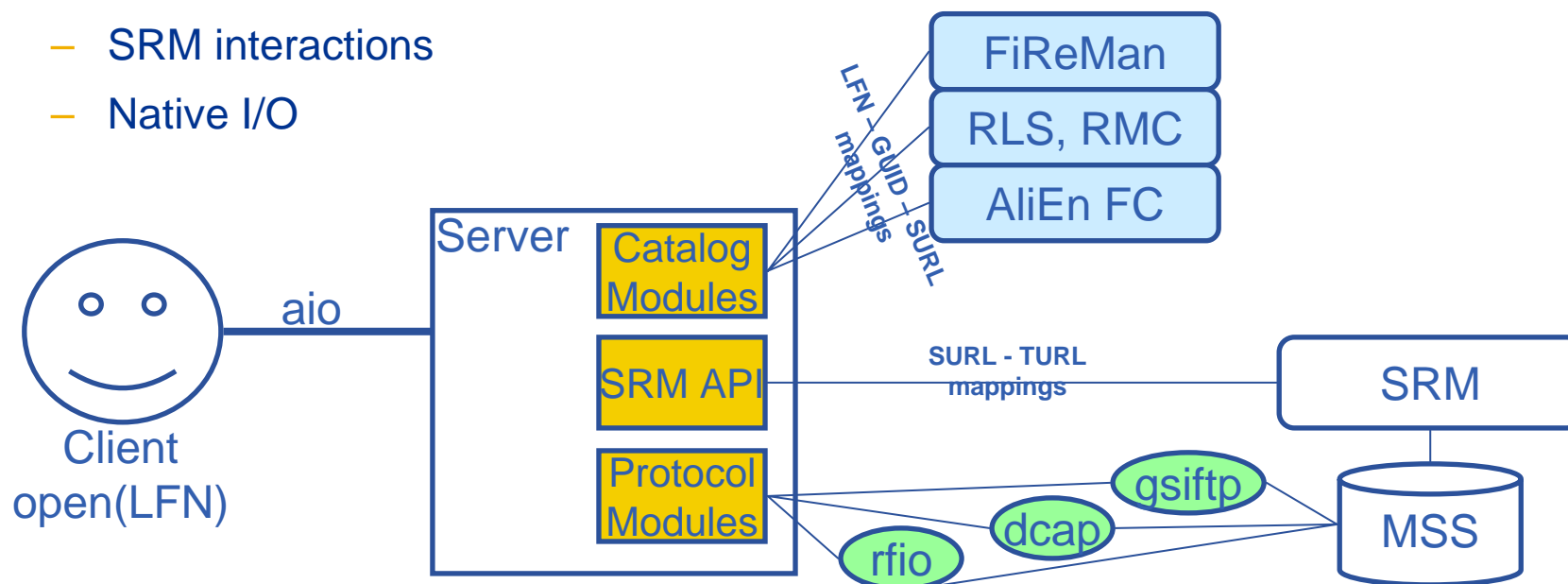
- StorageIndex: file location for broker
- FAS: File Access Service (ACLs)
- File Catalog: directory structure in LFN namespace
- Replica Catalog: location of replicas
- Meta: additional (user defined metadata)

Implemented on top of Oracle and MySQL



- **Grid Data Management Challenge**
- **Storage Elements (SRM)**
- **LFC (LCG File Catalog)**
- **File and Replica Catalog (FiReMan)**
- **gLite I/O**

- **Client only sees: a simple API library and a Command Line Interface.**
  - GUID or LFN can be used, i.e. `open("/grid/myFile")`
- **GSI Delegation to gLite I/O Server**
- **Server performs all operations on User's behalf**
  - Resolve LFN/GUID into SURL and TURL
- **Operations are pluggable:**
  - Catalog interactions
  - SRM interactions
  - Native I/O

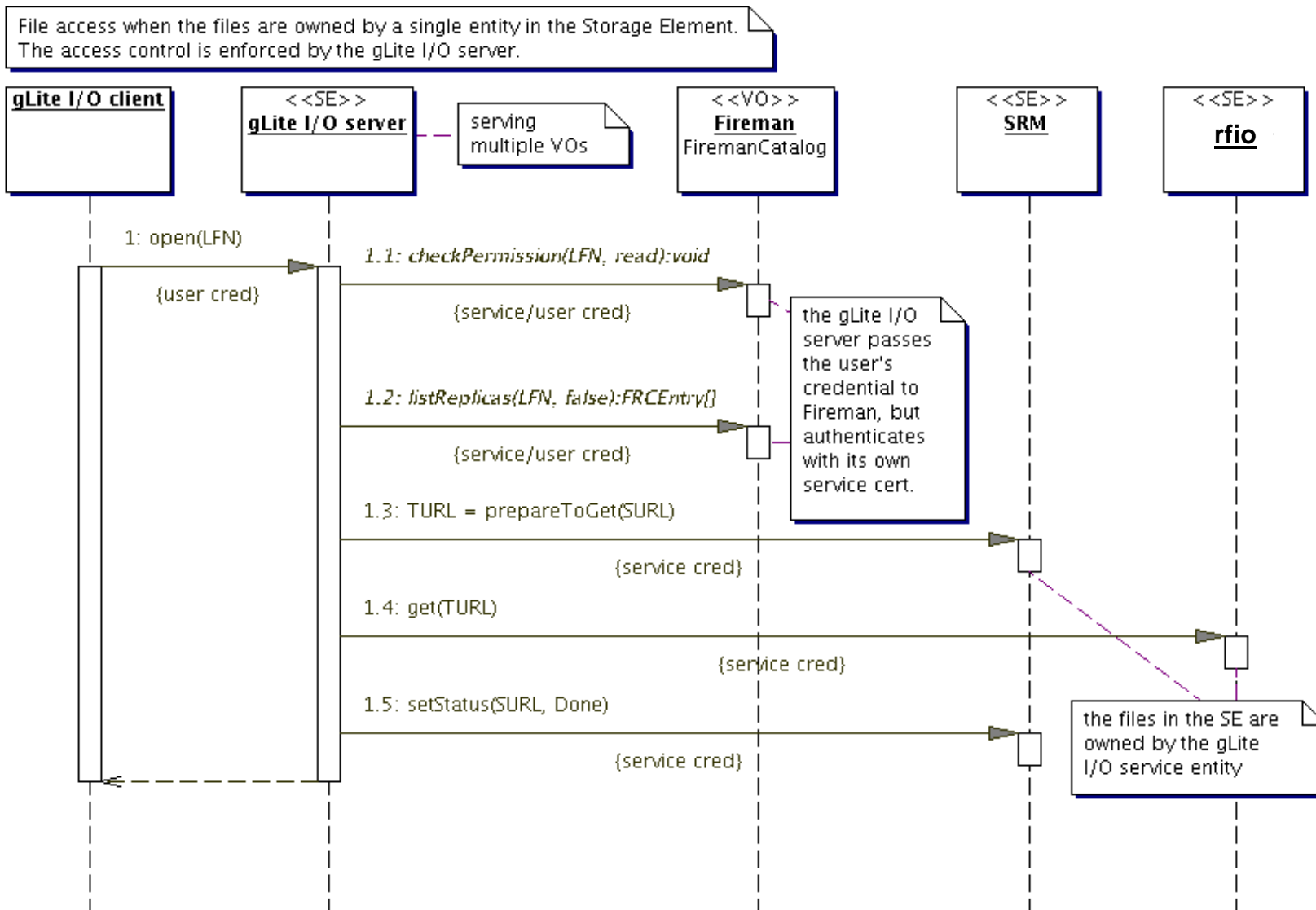


## Summary of the gLite I/O command line tools

<b>glite-get</b>	<b>Retrieve a file from the Grid using LFN or GUID</b>
<b>glite-put</b>	<b>Put a local file into the Grid, assigning LFN</b>
<b>glite-rm</b>	<b>Remove a file (replica!) from the Grid using LFN or GUID</b>

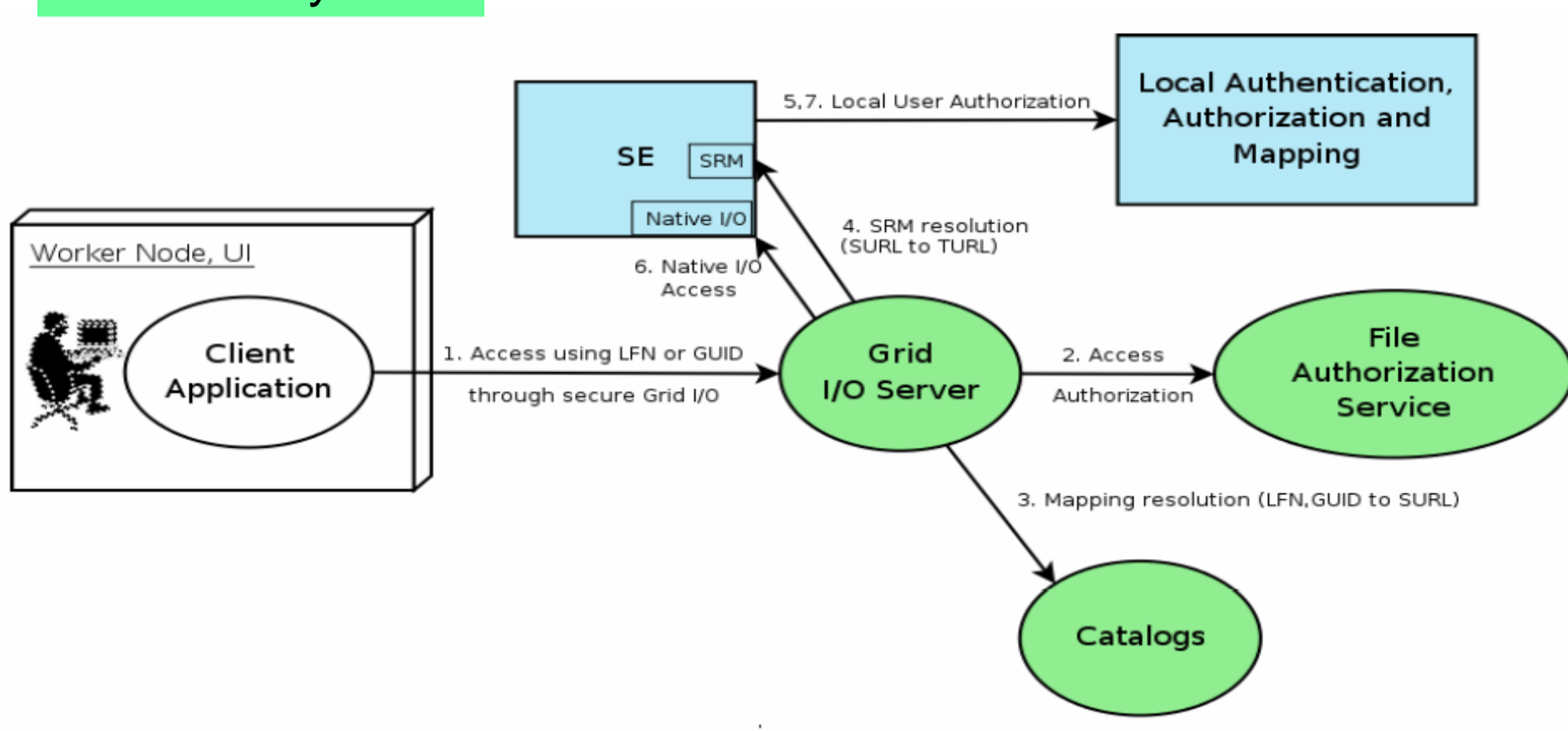
## Summary of the gLite I/O API calls (C only)

<b>glite_open</b>	<b>glite_posix_open</b>
<b>glite_read</b>	<b>glite_posix_read</b>
<b>glite_write</b>	<b>glite_posix_write</b>
<b>glite_creat</b>	<b>glite_posix_creat</b>
<b>glite_fstat</b>	<b>glite_posix_fstat</b>
<b>glite_lseek</b>	<b>glite_posix_lseek</b>
<b>glite_close</b>	<b>glite_posix_close</b>
<b>glite_unlink</b>	<b>glite_posix_unlink</b>
<b>glite_error</b>	<b>glite_filehandle</b>
<b>glite_strerror</b>	

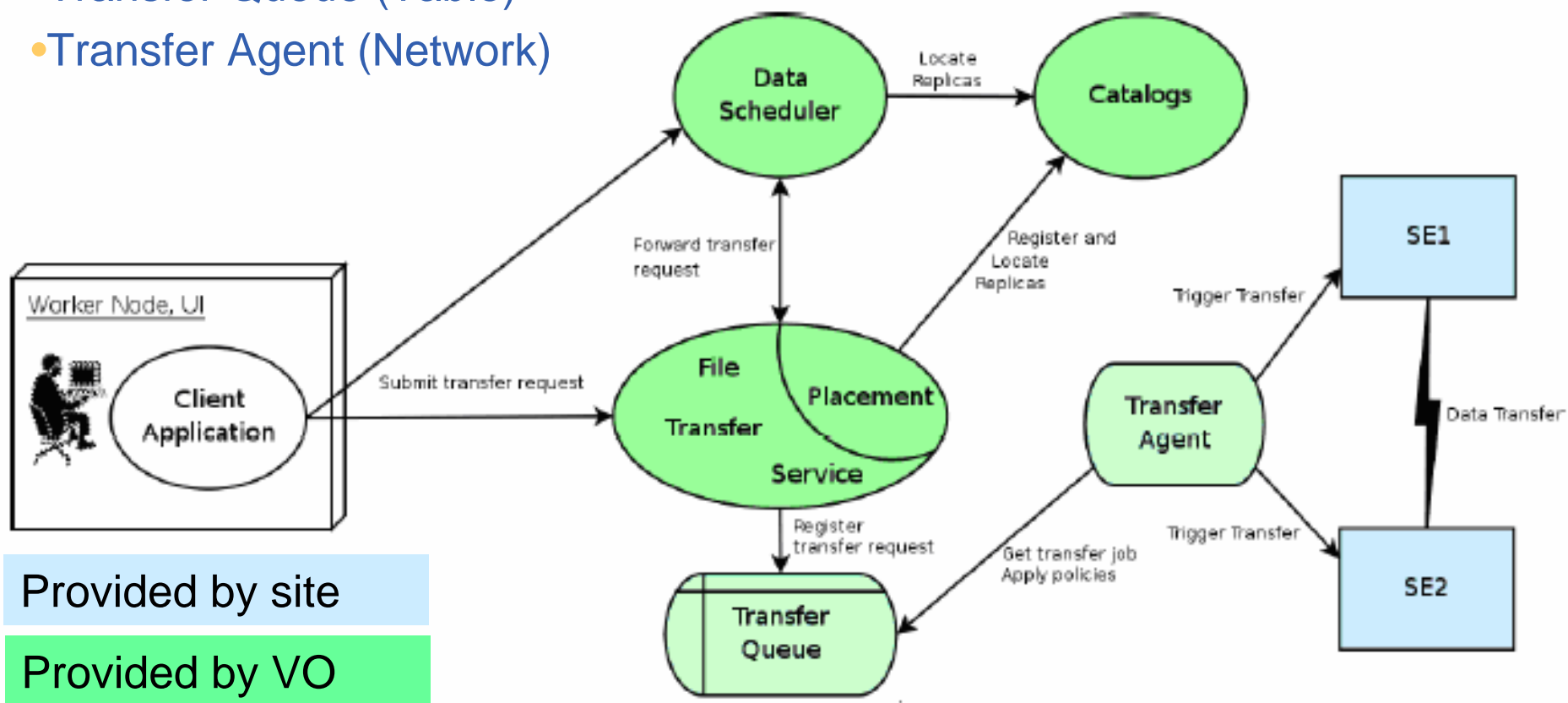


Provided by site

Provided by VO



- Data Scheduler (DS) Keep track of user/service transfer requests
- File Transfer/Placement Service (FTS/FPS)
- Transfer Queue (Table)
- Transfer Agent (Network)



- **gLite homepage**
  - <http://www.glite.org>
- **DM subsystem documentation**
  - <http://egee-jra1-dm.web.cern.ch/egee-jra1-dm/doc.htm>
- **FiReMan catalog user guide**
  - <https://edms.cern.ch/file/570780/1/EGEE-TECH-570780-v1.0.pdf>
- **gLite-I/O user guide**
  - <https://edms.cern.ch/file/570771/1.1/EGEE-TECH-570771-v1.1.pdf>

