



# ARDA Experiences with Metadata

Birger Koblitiz

LCG Application Area Meeting, January 26<sup>th</sup>, 2004

## Overview

- Metadata on the GRID
- Experiences with Experiment Solutions
- A Generic Interface to Metadata
- An ARDA Metadata Service Prototype
- SOAP: A good way to Access Metadata?
- Ideas for the Future
- Conclusions



# What is Metadata?

## 1. Definition:

Metadata is information on contents of files.  
(File Metadata)

Also other information found in DBs necessary to run jobs on the grid, share problems:

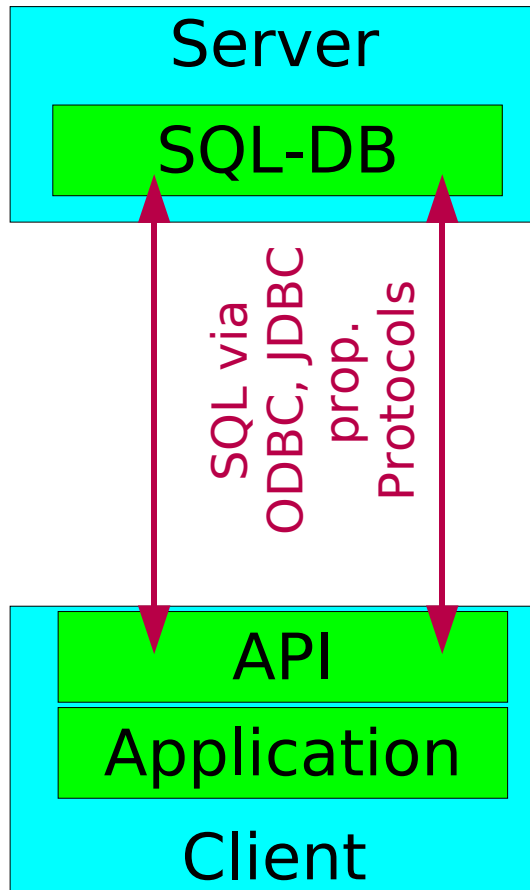
- Grid authentication
- Overcoming firewalls
- Talking efficiently to DBs
- Replication, different types of storage

## 2. Definition:

Metadata is data living in databases needed by jobs to run on the grid.

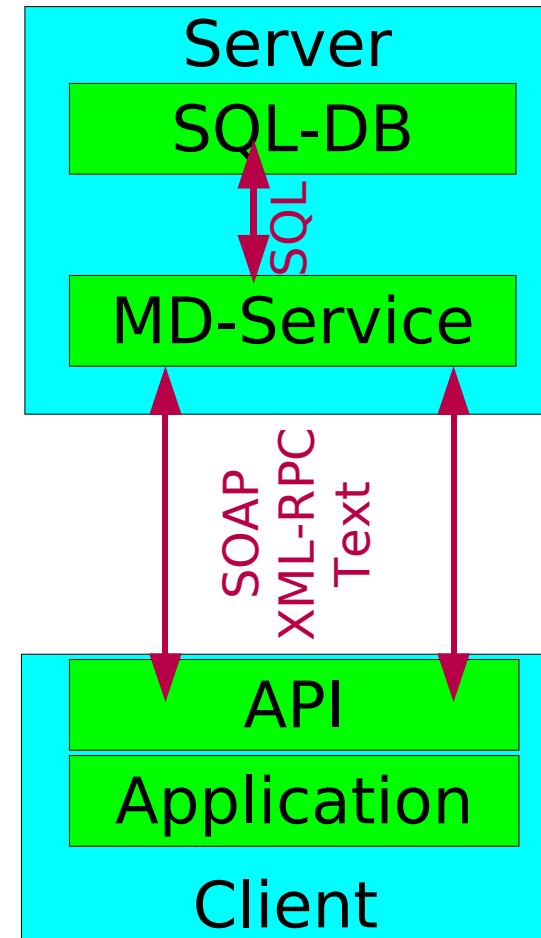
# DB Access on the Grid

“Traditional” Way: ODBC, RAL, ...



- +Performance
- +Simple Implementation
- Security, Monitoring

“Service”: AMI, RefDB, ...



- +Security: GSI, x509
- +Lightweight Client
- Performance, Impl.



# Experience

ARDA tested several Metadata solutions from the experiments:

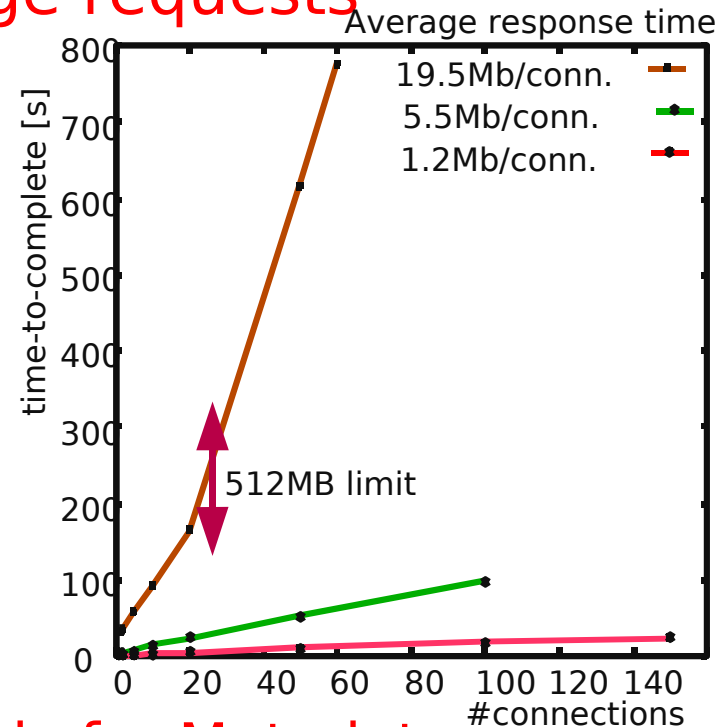
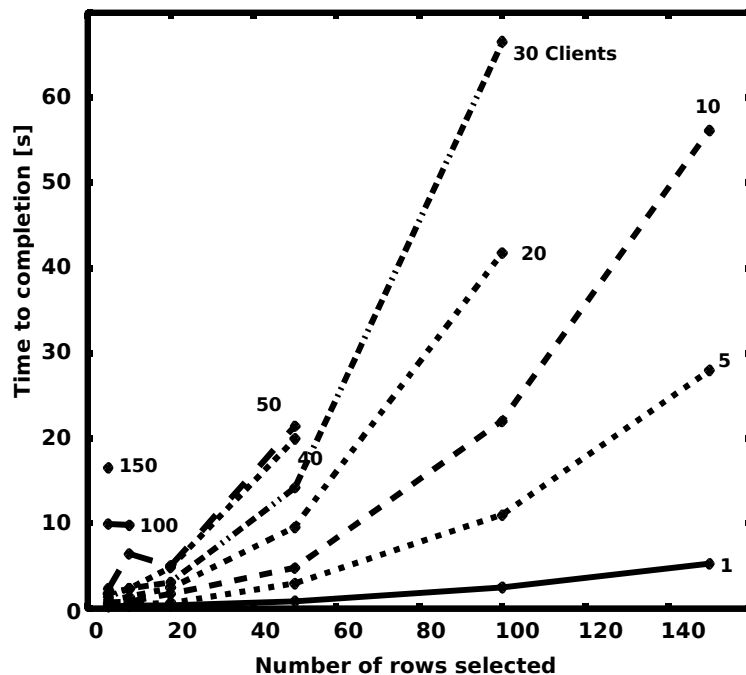
- LHCb Bookkeeping: XML-RPC with Oracle backend
- CMS: RefDB  
(PHP in front of MySQL, giving back XML tables)
- Atlas: AMI  
(SOAP-Server in Java in front of MySQL)
- gLite (Alien Metadata)  
(Perl in front of MySQL parsing command, streaming back text)

Learned a lot looking at existing implementations:

- Common pattern seen
- Implementations also share the same problems

# AMI & RefDB (SOAP)

Both **AMI & RefDB** ship responses in single XML package  
 → They **can't handle large requests**



**SOAP is particularly problematic for Metadata:**

- SOAP blows up data by factors 5-10
- SOAP for single, small queries

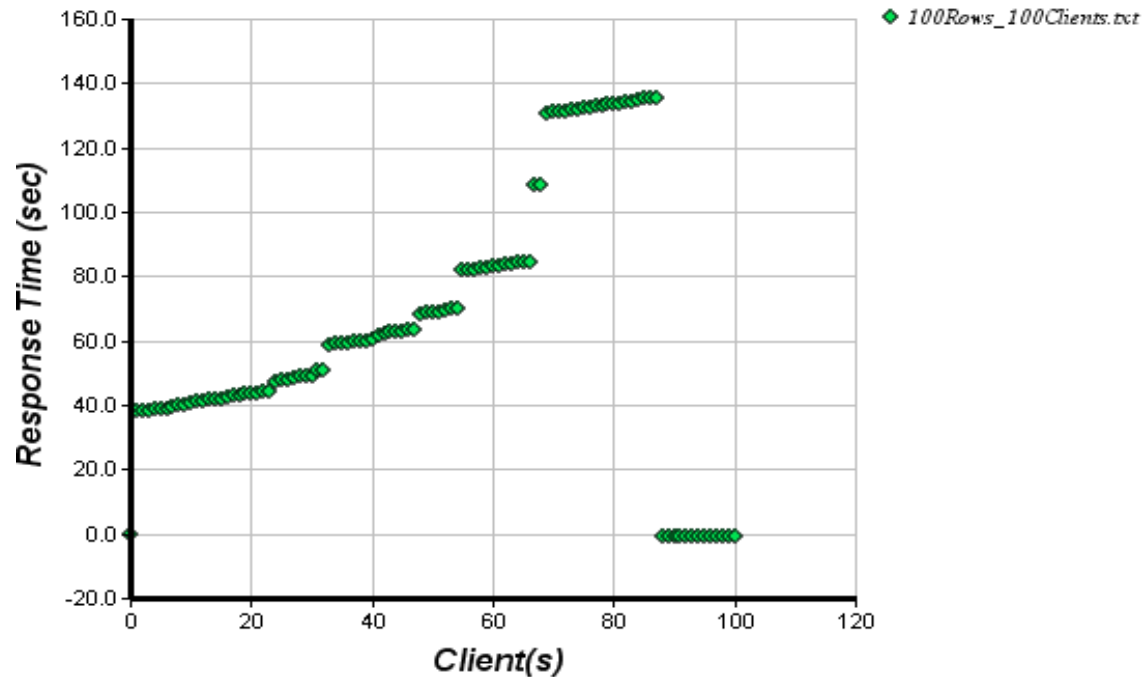
Metadata queries do require **stateful connections** with **Streamed Data / Iterators** as a response

Results on RefDB: J. Andreeva, J. Herralá

# LHCb(XML-RPC)

LHCb uses XML-RPC(predecessor of SOAP):

*LHCb Bookkeeping Testing Result*



created with ChartDirector from [www.advsofteng.com](http://www.advsofteng.com)

Being also “one shot” query based, the solutions suffers from the same problems as the two based on SOAP

# SOAP extreme

Snippet of client code from gLite Fireman:

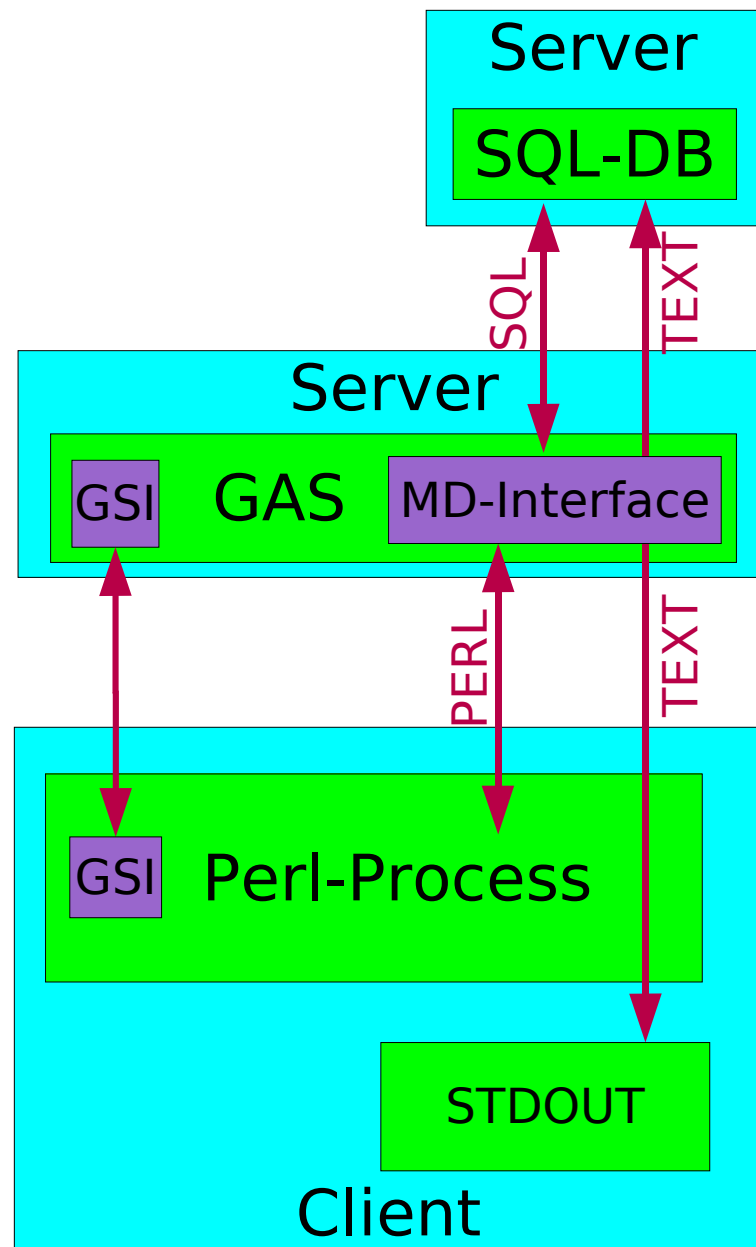
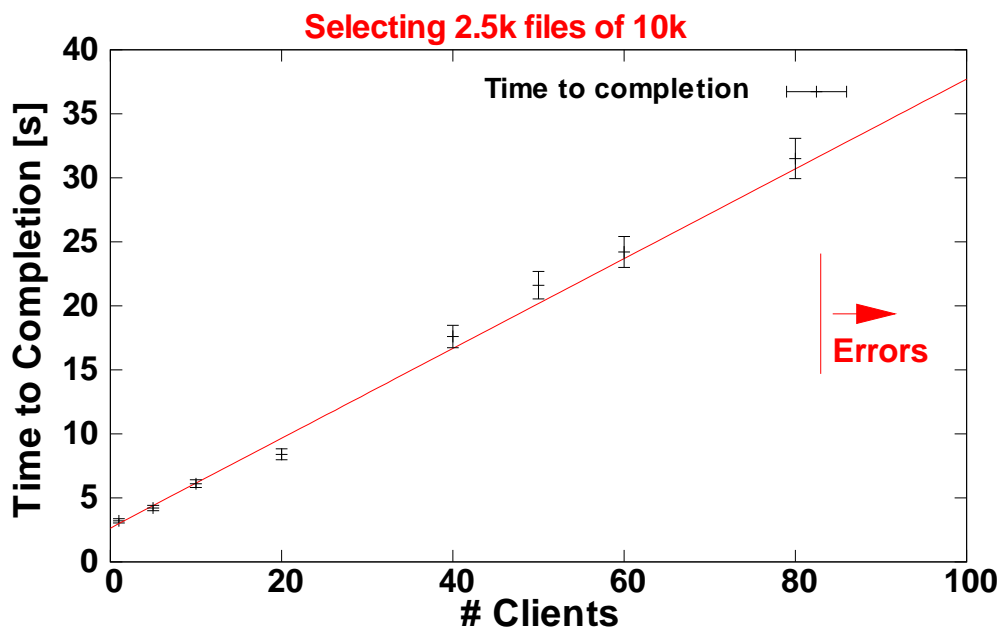
```
file::ArrayOfFCEntry fc_entry;
fc_entry.__size      = 1;
fc_entry.__ptr      = (file::glite_FCEntry**)
    soap_malloc(fileService.soap, sizeof(file::glite_FCEntry*));
fc_entry.__ptr[0] = file::soap_new_glite_FCEntry(fileService.soap, -1);
fc_entry.__ptr[0]->guid      = soap_strdup(fileService.soap, guid);
fc_entry.__ptr[0]->lfn      = soap_strdup(fileService.soap, lfn);
fc_entry.__ptr[0]->permission = 0;
fc_entry.__ptr[0]->lfnStat  = file::soap_new_glite_LFNStat(fileService.soap, -1);
fc_entry.__ptr[0]->lfnStat->type      = 0; // LFN
fc_entry.__ptr[0]->lfnStat->data      = 0; // Additional Information
fc_entry.__ptr[0]->lfnStat->modifyTime = 0; // Use Default Value
fc_entry.__ptr[0]->lfnStat->validityTime = 0; // Use Default Value
fc_entry.__ptr[0]->lfnStat->creationTime = 0; // Use Default Value
file::file__createFileResponse out;

if(SOAP_OK != fileService.file__createFile(&fc_entry, out)){
    // TODO Exception Handling
    // Finalize service Object
    // finiFileService(&fileService);
    return -1;
}
```

Complex (bulk ?) SOAP calls difficult to use without API!  
Several incompatibilities among SOAP implementations!

# gLite/Alien(Streaming)

gLite/Alien streams responses to the perl implemented shell







# Schema Handling

Schema evolution not really tackled by current Metadata Catalogues:

- Not really important for production...
- Admin can setup/copy new tables (work on backend)...

**RefDB** and **Alien** don't do schema evolution at all.

**AMI, LHCb**-Bookkeeping via admins adjusting tables.

**EGEE design for gLite** does not foresee schema changes nor schema discovery

For analysis, the following capabilities are mandatory:

- User must be able to **discover schema**
- User can **setup/change schema**
- Offer **solution for problems with storage types**



# POSIX Metadata

POSIX defines **extended attributes (Metadata)** for files:

- **Key-Value pairs** associated with a file
  - Key: \0-terminated string
  - Value: Binary data of arbitrary length
- Copying a file copies metadata
- Metadata can be attached to directories (no inheritance)
- Metadata attached to inode (security)

Extended attributes are now widely used  
(NTFS, NFS, EXT2/3 SCL3, ReiserFS, JFS, XFS)

Used with Namespaces for ACLs

Metadata searches not defined yet (No FS-Impl.):

- Windows Longhorn (2005)
- ReiserFS 5



# Metadata on Linux

On ext3, XFS or ReiserFS, Linux supports extended attributes (file metadata)

```
koblitz@pcardabk:~/test$ touch a
koblitz@pcardabk:~/test$ attr --help
Usage: attr [-LRSq] -s attrname [-V attrvalue] pathname # set value
        attr [-LRSq] -g attrname pathname             # get value
        attr [-LRSq] -r attrname pathname             # remove attr
        -s reads a value from stdin and -g writes a value to stdout
koblitz@pcardabk:~/test$ attr -s gen -V lept0 a
Attribute "gen" set to a 5 byte value for a:
lept0
koblitz@pcardabk:~/test$ attr -s version -V 1.0 a
Attribute "version" set to a 3 byte value for a:
1.0
koblitz@pcardabk:~/test$ getfattr -d a
# file: a
user.gen="lept0"
user.version="1.0"
koblitz@pcardabk:~/test$ grep home /etc/fstab
/dev/hda5 /home ext3 defaults,acl,user_xattr,auto 0 0
```

Can we have a similar semantics on the Grid?

PS: API is POSIX, not the commands!

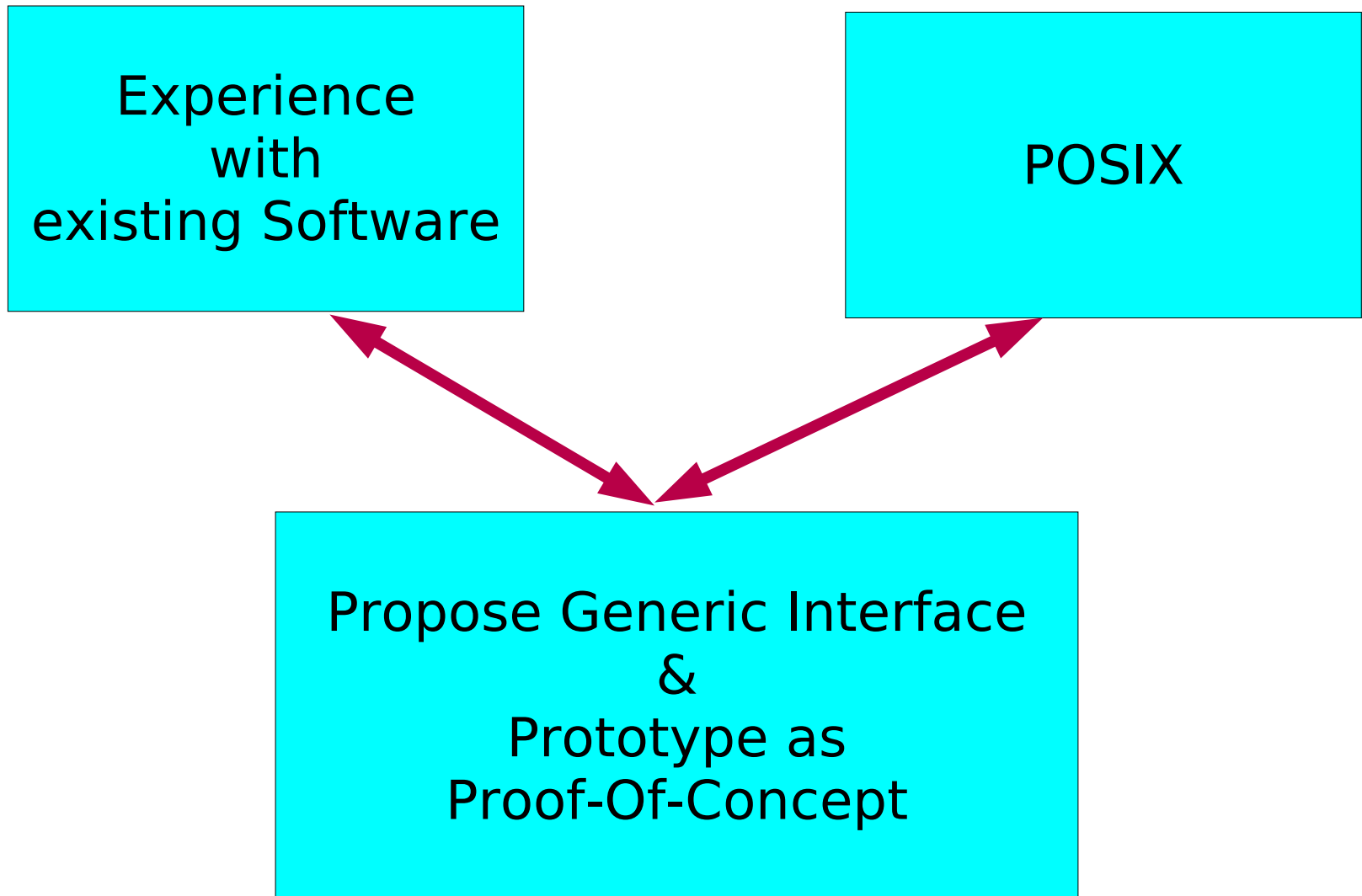


# GRID Metadata

A possible Grid approach:

- Metadata attached to LFN
  - LFN is entry point to File-Catalogue, attached Metadata can be easily searched
- Files without LFN: GUID in special dirs
  - Otherwise problems with global searches
- Metadata for directories should provide default schemas/values for files
  - Easy schema copying
- Restrict values to ASCII strings
  - Backend is unknown: FileSystem/DB
- Need to define ways how to search for Metadata: Search restricted to (sub-)directories
  - Allows hierarchical databases, applicable for FS

# Synthesis





# Hierarchy

Metadata needs a hierarchy to work well:

- Collect objects with shared attributes into collections
  - ➔ Allows queries on SQL tables  
(also other storage possible: XML-DB, DB-Files...)

Analogy to file system (file metadata!):

Collection ↔ Directory

Object ↔ File

- Structure important for:
  - Structured searches
  - Schema handling
  - Distribution of databases



# Terminology

Define common terms, first

- **Metadata: Key-Value pairs**  
Any data necessary to work on the grid not living in the files
- **Entry: Entities to which metadata is attached**  
Denoted by a string, format like file-path in Unix, wild-cards allowed
- **Collection: Set of entries**  
Collections are themselves entries, think of Directories
- **Attribute: Name or key of piece of metadata**  
Alphanumerical string with starting letter
- **Value: Value of an entry's attribute**  
Printable ASCII string
- **Schema: Set of attributes of an entry**  
Classifies types of entries: Collection's schema inherited its entries
- **Storage Type: How back end stores a value**  
Back end may have different (SQL) datatypes than application
- **Transfer Type: How values are transferred**  
Values are transported as printable ASCII



# Interface I: Entries

Interface designed in discussions with gLite-team, GridPP metadata working group, GAG, submitted to PTF

The following protocol is proposed which clients talk to servers via sockets:

- **int addEntry(string entry, string type)**

Adds a new **entry** to the catalogue

**Type** can be “Collection” or “Entry”

(extensions from implementation: Inheriting collections, views...)

Returns integer errors code: MD\_SUCCESS=0, MD\_ERR\_NOENT, MD\_ERR\_PERM, MD\_ERR\_INT

- **int addEntries(list<string> entries, list<string> types)**

**Entries** and **types** lists for bulk insertion into catalogue

- **int removeEntries(string pattern)**

**Pattern** for intuitive bulk deletion

With R. Rocha(gLite), V. Pose, N. Santos





# Interface II: Attributes

Schema management and metadata reading/writing:

- **int addAttr(string entry, list<string> keys, list<string> types)**

Adds a new attribute (**key**) to an **entry** (collection)

- **int removeAttr(string entry, list<string> keys)**

Removes attributes (**keys**) of an **entry** (collection)

- **int setAttr(string pattern, list<string> keys, list<string> values)**

Bulk setting of the **keys** of all entries matching **pattern** to new **values**.

- **int clearAttr(string pattern, string key)**

Resets the **keys** of all entries matching the **pattern**  
Application will get empty string if asking for value.  
Behaviour in queries like NULL in SQL.

# Interface: Retrieving data

The Bulk transfers to client are done through **session handlers** and **iterators** on the backend:

- **int getAttr(string pattern, list<string>keys, Handler &handler)**

Returns values for all **keys** of the entries matching **pattern**

```
struct Handler {  
    handle_t handle;  
    DataChunk chunk;  
}  
  
struct DataChunk {  
    list<string> values;  
    bool last;  
}
```

The values contain names of matching entries end the data:

→ Client knows semantic

- **int getNext(handle\_t handle, DataChunk &c)**

Returns the next bunch of values

- **int abort(handle\_t handle)**

Aborts request

- **int listAttr(string entry, Handler &handler)**

Lists all attributes of an entry

With R. Rocha(gLite), V. Pose, N. Santos

# Interface: Searching

Physics analysis needs powerful tool to find entries, more than attribute-value matching:

- **int find(string pattern, string query, Handler &handle)**

Returns all entries (excluding collections) matching the **pattern** and fulfilling the **query**

```
Example query: 'tracks > 10 and sin(p_angle) <0.5 and trigger & 2'
```

Query needs to be parsed:

- SQL injection prevention
- Separate user & system namespace: events → “user:events”
- Interpret for different backends

- **int listEntries(string pattern, Handler &handle)**

Returns all entries and collections matching the **pattern**, giving their type



# Protocol Choices

Presented interface is **SOAP compatible**:

- Fulfil formal requirements on EGEE
- SOAP implementation using gSOAP for server done:  
Clients exist in C++, Python, Java

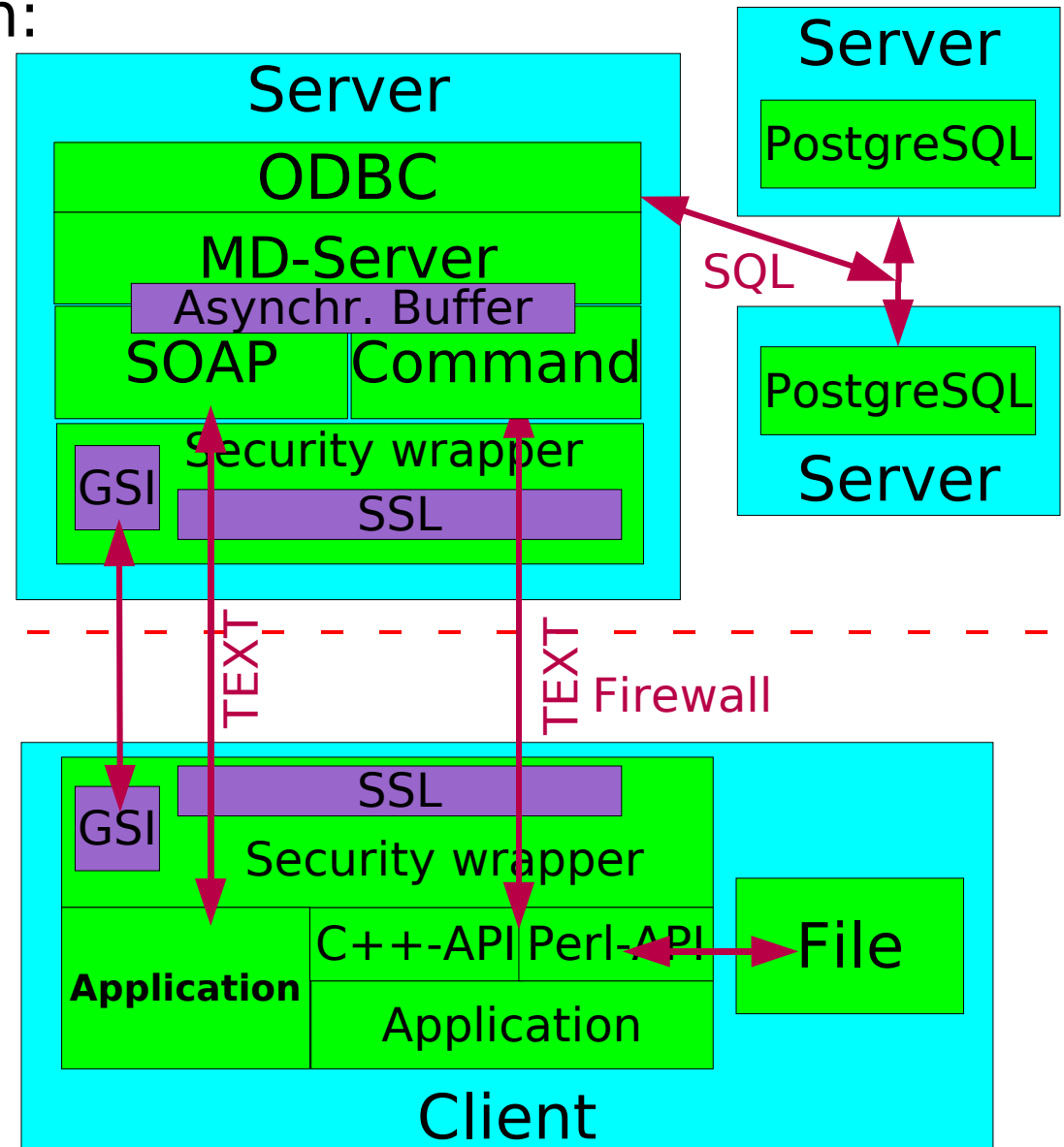
Also implemented as TCP/IP streaming:

- Use plain text (**ASCII**)
  - Query consists of **one line of command**
  - Response returns 1 line of return status (OK/Error) and result line by line (and EOT at end)
  - Result is in ASCII, user needs to encode/decode
    - Commands are Line of ASCII, e.g:
      - `getattr file(s) key1 key2...` Returns value of keys
  - SSL for authentication and encryption implemented
  - Clients in C++, Java, Python, Perl, Ruby
- Implementing client APIs very simple!**

# Prototype

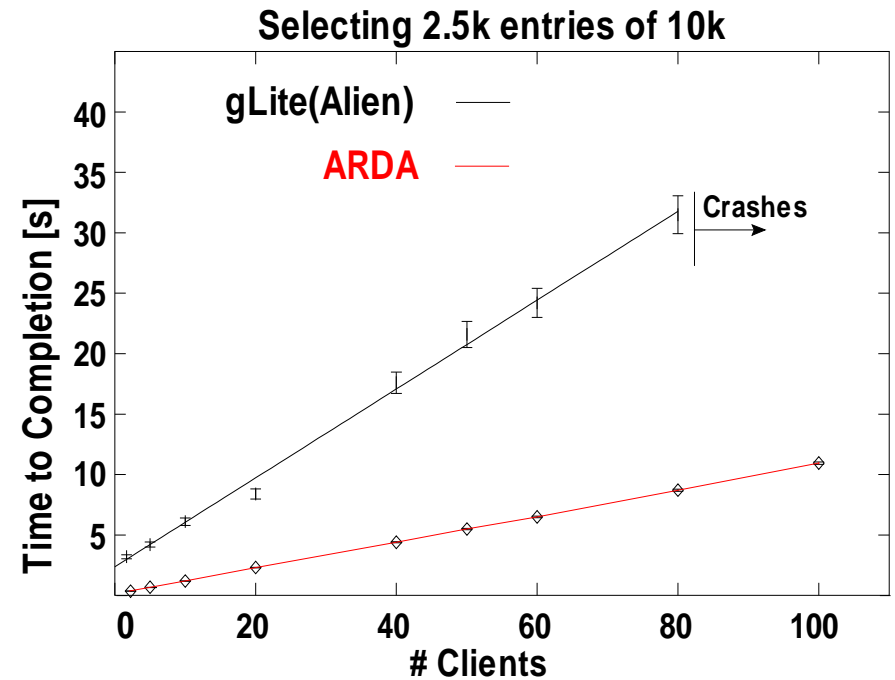
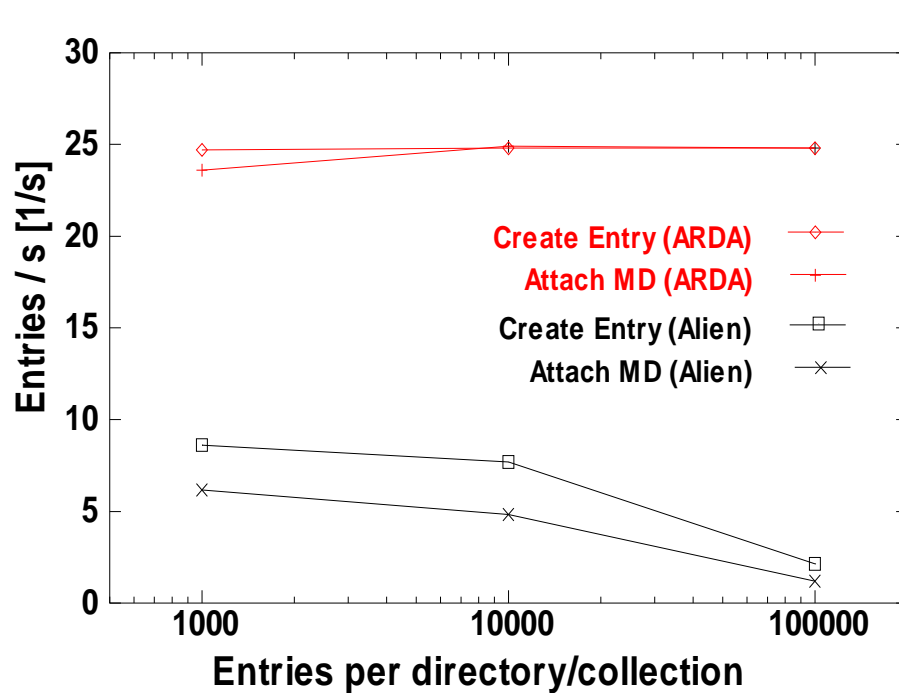
## Prototype Implementation:

- Multi-threaded C++ server in front of PostgreSQL
- Streams responses asynchronously
- Uses ODBC as RDBMS abstraction Layer: ODBC-types
- Access restrictions via ACLs
- Bison/flex parser for queries
  - Other backends
  - Query validation
  - Security



# Reality Check

Good experiences with ARDA prototype using streaming:



ARDA prototype even faster than AliEn!

Now very stable after tuning.

Started Collaboration with LHCb, large scale test till

January 2005

New GANGA back end?

With V.Pose



# SOAP Test

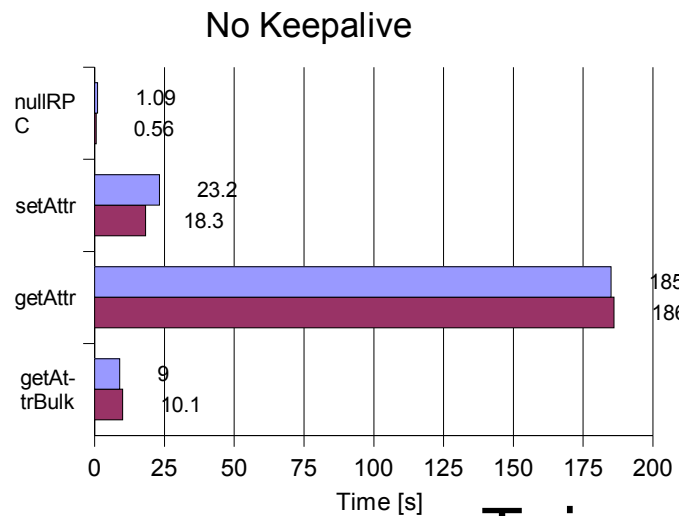
Have sophisticated Metadata prototype with

- SOAP+Streaming front ends
- Identical back ends

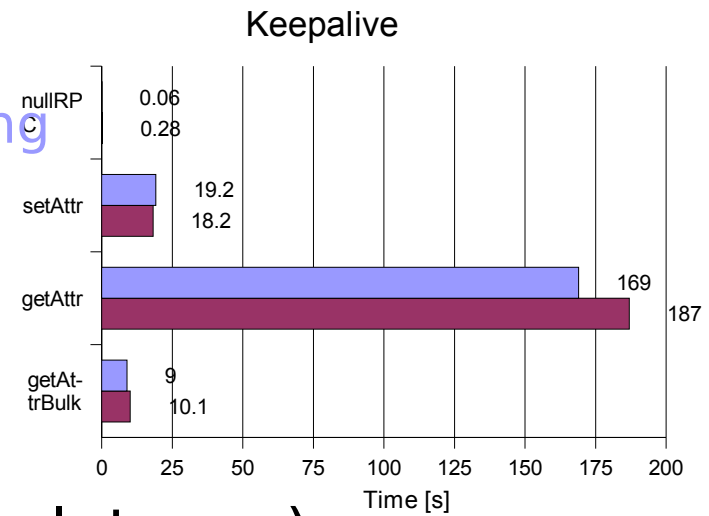
→ Study SOAP as metadata access protocol

# SOAP vs Streaming

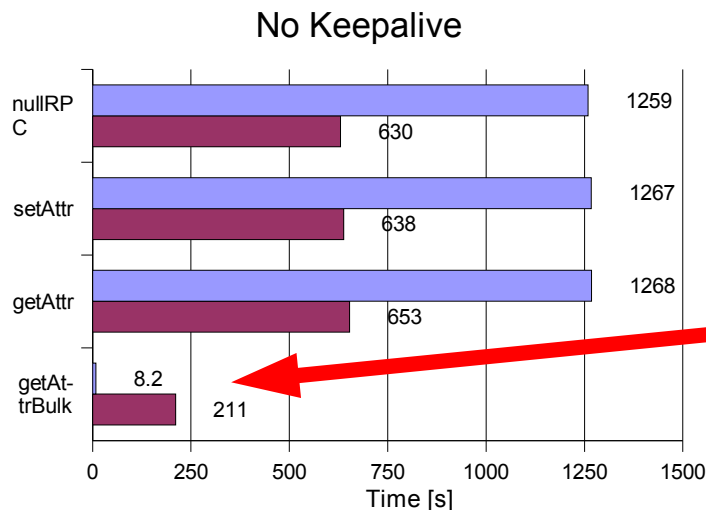
C++-Client (gSOAP): Create/read 1000 entries/60 attrs.  
Locally (same computer)



Streaming  
SOAP

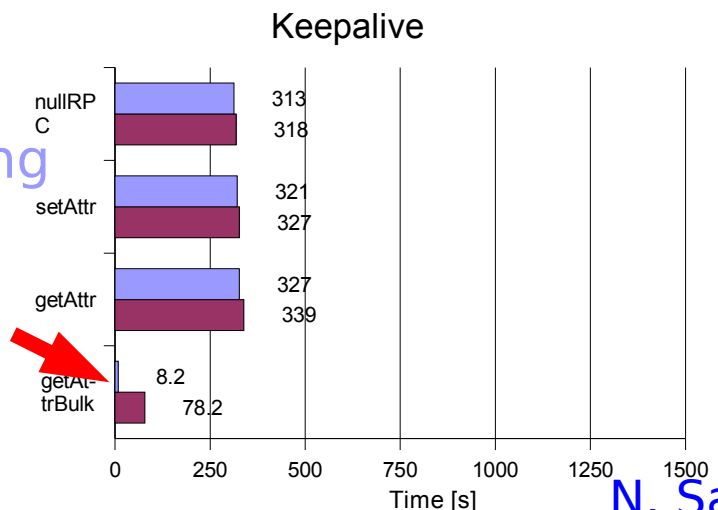


Taiwan (312ms latency)



Streaming  
SOAP

x25!



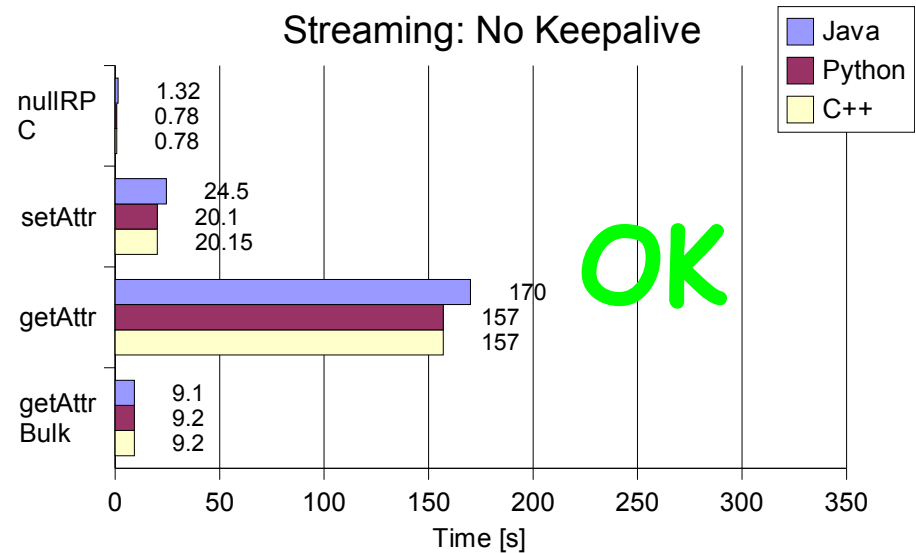


# SOAP Toolkits

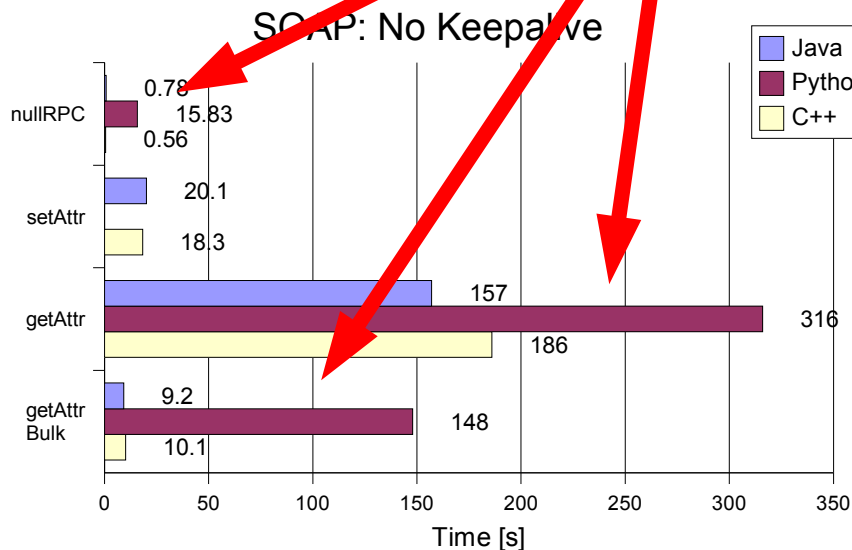
Create/read 1000 entries/60 attrs. locally

No performance differences of Socket implementations.

SOAP toolkits immature?



OK



Took experienced programmer week to create interoperable WSDL:

Complex datatypes problematic.



# Conclusions on SOAP

Preliminary tests show good performance of SOAP in single shot remote calls.

Retrieving large results is problematic.

SOAP toolkits not all mature, but gSoap highly optimized.

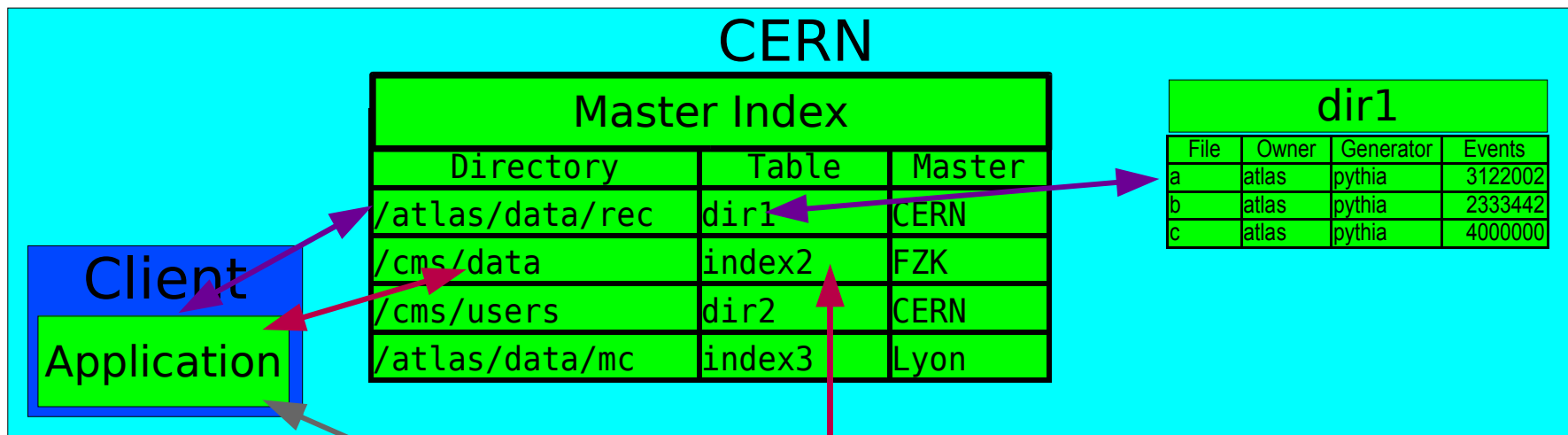
Need to study servers under load by many clients: Observed factor of 6 in increased traffic.

SOAP interoperability very problematic!

# Distributed Metadata Ideas

Currently lots of brainstorming done:

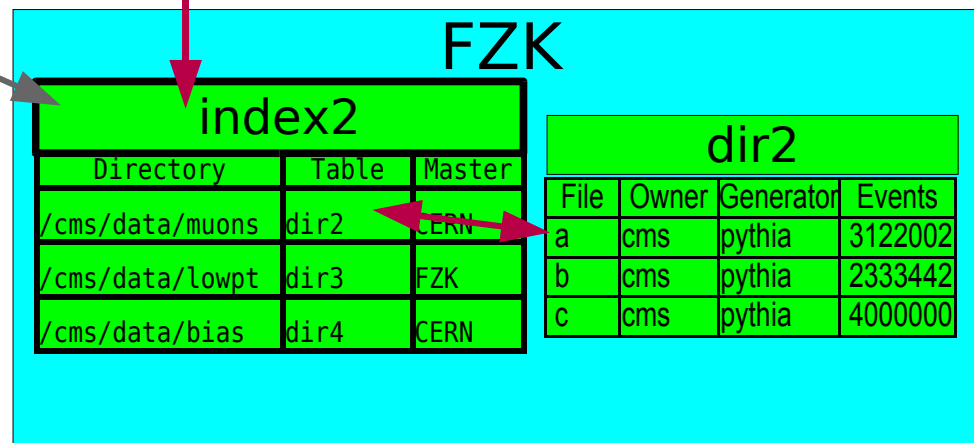
Use PostgreSQL/Oracle **per-table replication with different masters**,  
make use of hierarchy:



Talk directly?

Thesis: N. Santos  
Steering possible through  
common interface...

Participating in LCG 3D.





# Conclusions

- Many problems understood studying metadata implementations of experiments
  - Common requirements exist
  - ARDA proposes generic interface to metadata on Grid:
    - Retrieving/Updating of data
    - Hierarchical view
    - Schema discovery and management  
(discussed with gLite, GridPP, GAG, submitted to PTF)
  - Prototype with SOAP & streaming front end built
    - SOAP can be as fast as streaming in many cases
    - SOAP toolkits need time to mature
  - Efficient DB access on Grid challenging
    - Problems: Distribution, Authentication, Monitoring,...
    - Metadata experts need to work together
- <http://project-arda-dev.web.cern.ch/project-arda-dev/metadata/>