Modeling Storage Resources with GLUE information providers

Jeff Templon, Pierre Girard, Stephen Burke, Peter Kunszt

January 10, 2006

1 Introduction

This document is motivated by problems observed in trying to accomplish two use cases on grid-enabled storage resources. While trying to figure this out, it became clear that

- the current software stack is using only a small fraction of the GLUE information published,
- most of the resources are publishing only a fraction of the GLUE attributes
- the information that is published is often at variance with the GLUE specification, and
- it is not at all clear how one is supposed to model a real SRM-type SE using the GLUE schema.

A good example of this last point is the GLUE **StorageArea** class, which seems to mix virtual concepts like quotas per VO and SURL paths, with physical concepts such as capabilities to pin files onto a disk cache area.

We hope that this document will not only result in a proposal on how to accomplish the use cases, but also in a description of how one should actually model a real storage element using the GLUE schema, and inform the GLUE development process for the next round.

1.1 How to find disk SEs

LHCb wanted to know how to find the space at each site that was 'pure disk', meaning that if they wrote a file there, it would stay on disk and never need to be staged out from tape. Evidently this was not easy via grid methods since the suggested solution was to define two SEs,

- se-disk.nikhef.nl and
- se-mss.nikhef.nl,

effectively using the SE hostname as storage element metadata.

1.2 How to ensure pinning (in cache area) of files

Several groups have expressed interest in having files, stored in MSS-based systems, being 'pinned' on fast cache areas for long times. It's not clear at all how to do this without contacting a center and specifically requesting the staff to do this.

2 Assumptions

The goal of GLUE seems to be, at the highest level, a description of a virtualized storage resource, independent of the underlying hardware. We try to align our proposal with this goal as long as it remains practical to do so. This also agrees with the goal of the SRM interface, which is supposed to be largely virtual (independent of the underyling storage architecture).

3 Definitions

Reports are that the SRM team put a lot of thought into terminology, so we propose to use that if it really exists. Not all authors have seen the 'definitive SRM nomenclature document' yet but we'll give its existence the benefit of the doubt here.

Please try to use, in these discussions, the terms below only in the sense described below. Using them in other ways is a quite effective technique for creating counterproductive discussions.

- **SURL** Storage URL. We use this here to mean the URL itself, as well as the file it refers to. The existence of a SURL in a given SE just means that it is possible to access that file from that SE (interaction with other SEs is not needed, neither by the user nor the SE in question). However it says nothing about whether the file is permanently on disk, on both back-end slow archival storage and on disk cache space, or only present on the back-end archival store.
- **TURL** Transfer URL. We use this here to mean the URL itself, as well as the file to which it refers. A TURL is a file on an SE that can be accessed immediately. This probably means that it is either located on a disk-based SE (all files are permanently active) or else located on the cache area of a MSS system.
- **online** this means that the file is currently located on fast storage, hence one could receive a TURL for it and access it immediately without waiting for it to be staged in from archival medium.
- online staging area We've seen 'cache' and 'buffer' used for this term as well, online staging area is SRM terminology. We *think* It refers to the area where the SE creates TURLs when it stages them in from archival medium. SURLs which have corresponding TURLs in the online staging area are online in the sense defined above.
- **nearline storage** the archival medium of a MSS. The implication is that this is slow storage, ie it might take time to get it out into the online staging area and return a TURL.

4 Disk-Based vs. MSS-Based SEs

There is an easy way to distinguish between the two: the Architecture property of the GLUE StorageElement class. Allowed values listed in the GLUE 1.2 specification are 'disk, tape, multidisk, other'. 'multidisk' is not being used on LCG at the moment; it refers to disk-based systems that are operating in a raid-like redundant configuration.

It is not clear that looking for Architecture type 'disk' or 'multidisk' is what an experiment really wants to do; what they really want in this case is to put the data somewhere so that it will be quickly accessible for a long time to come. There may be disk systems for which this is not true, *e.g.* a DPM-based disk SE with a single network interface and a limited number of allowed incoming transfer connections; you might have to wait before you are given a TURL.

5 What to do about online staging areas

Both SRM and GLUE SE assume they are trying to come up with a virtualized description of a storage resource. Given this goal and the current organization of the GLUE SE schema, it seems we have to conclude that the GLUE **StorageArea** refers to, as it is now (January 2006), a space for SURLs; it is the only part of the SE schema that refers to VOs at all.

A related question, discussed at length during EDG but never resolved, is whether a SURL should uniquely identify a file. If SURL paths look like

```
srm://dm.lbl.gov/my.file
```

then we might have six of these critters, each assigned to a different VO. Things get tricky if we ever need to have VOs access each others' files. Not important for HEP most likely, but probably important for many other groups. On the other hand,

```
srm://dm.lbl.gov/atlas/my.file
```

would work since the VO name is part of the path. Current LCG usage is inconsistent. Here is a list of the various types of (GlueSAPath) values currently declared for ATLAS:

```
GlueSAPath: atlas:/atlas
GlueSAPath: atlas:/castor/grid.sinica.edu.tw/sc/atlas
GlueSAPath: atlas:/pnfs/gridpp.rl.ac.uk/data/atlas
GlueSAPath: atlas:/pnfs/gridpp.rl.ac.uk/tape/atlas
GlueSAPath: atlas:/pnfs/usatlas.bnl.gov/data
GlueSAPath: /castor/ific.uv.es/grid/atlas
GlueSAPath: /dpm/scotgrid.ac.uk/home/atlas
GlueSAPath: /grid/atlas
GlueSAPath: /grid/fzk.de/mounts/nfs/data/lcg1/SE00/atlas
GlueSAPath: /hpss/in2p3.fr/grid/atlas
GlueSAPath: /pnfs/grid.sara.nl/data/atlas
GlueSAPath: atlas
GlueSAPath: atlas:/dpm/itep.ru/home/atlas
GlueSAPath: /castor/pic.es/grid/atlas
GlueSAPath: atlas:data2/atlas/
GlueSAPath: /storage/atlas
```

Note the inconsistent usage of a) specifying the VO name in the path part, b) leading '/' characters on the path, c) specification of the site name in the path, and d) including the "voname: " prefix. We should decide on what the correct behavior is, now!

We suggest the following: the GlueSAPath is a component in how a program should try to access the file. We expect to be able to use it like this:

srm://GlueSEUniqueID/GlueSAPath/users_file

for a SRM based SE, and

gsiftp://GlueSEUniqueID/GlueSAPath/users_file

for a gridftp-based "classic" SE.

5.1 Minimal Solution for Second Use Case

One possible minimal extension is to add a new GlueSAState field: GlueSAStateOnline (boolean). If set to TRUE, this would mean that all the SURLs stored in the GlueSAPath corresponding to this StorageArea are online. This would satisfy the second use case: write to MSS but in a place that guarantees the files will remain on online staging storage as well. The disadvantage of this is that it is a real GLUE modification, and in principle we can't put this into production before going through the GLUE modification process.

The SARA operations team has said that dCache has this sort of operation: "files in this directory will not leave the disk". It would be good to know whether this is dCache specific behavior. How does e.g. CASTOR work?

5.2 General Concerns

The solution above costs little effort but only works for what we expect is a small subset of the usage of an SE. In general an SE with back-end nearline storage will have various available quotas for VO/groups on the nearline store(s) (GlueSAs) as well as various available quotas and policies for online staging areas. As far as we can tell there are more questions than answers here.

- Common usage: for SRMs in production now, is there a one-to-one mapping between SAs and online staging areas? Is there such a thing as a 'staging space' (chunk of online storage dedicated to a single VO, like a physical disk partition) or is it more like a GlueCE that supports several VOs??
- Do we limit the online staging area concept to where it makes sense (only in the case of systems with nearline back-ends) or do we insist on uniformity and have 'fake online staging areas' for disk-based SEs? It seems more logical to have disk SEs with no staging object defined.
- How do we model an OSA? Do we define another StorageArea with a 'onlineStagingArea' or 'cache' type? Or do we define a Glue CacheArea or OnlineStagingArea subclass of the SE schema?

- What do 'pin durations' mean for disk-based SEs? And what do they mean for the other type of SE? Does it mean how long the archival storage is guaranteed? Or does it mean for a surf stored in this SA, how long you are allowed to pin it as a turf in the corresponding CA?
- for the GlueSAType parameter, the specification says "guarantee on the lifetime for the storage area". Does this really mean this? If NIKHEF has a StorageArea cmsflammable, declared with GlueSAType: volatile then is NIKHEF free to delete the entire StorageArea at any time? Or does it refer to files stored under this area? Probably what a user wants to know is "under which area can I put files if I want them to have a certain lifetime property". Not "which StorageArea might completely disappear at any time".

6 Various

The stuff below has yet to find a place.

Assumptions:

pinning only refers to files in CACHE space, not the files "put" into the SRM. from notes from

http://sdm.lbl.gov/srm-wg/meeting0509/SRM-wg-050914.html

6.1 Definitions of GlueSAType parameter values

from http://www.nesc.ac.uk/events/GGF10-DA/programme/papers/SRMInterfaceSpecification.pdf

Note: "items" below can be files or spaces; not exactly clear about spaces, since the first doc above (sdm.lbl.gov link) has conflicting statements about whether spaces have types or not; clearly GLUE thinks that they DO have types since there is an SAType parameter.

permanent can only be removed by user/VO owning the item

volatile can be removed when its lifetime expires

durable can in principle be removed when lifetime expires, but one must take some action is either notifying the owners that removal is imminent, or first making an archive copy before removal.

If the above three types can refer to files and spaces, one needs to make clear what the implications are; it is also unclear whether the 'files' on cache space of an SRM have the same types (or types at all) as do the corresponding files on the back end system.