



NORDUGRID-TECH-2

3/12/2008

THE NORDUGRID GRID MANAGER AND GRIDFTP SERVER

Description and Administrator's Manual

A.Konstantinov*

*aleks@fys.uio.no

Contents

1	Introduction	3
2	Main concepts	3
3	Input/output data	4
4	Job flow	4
5	URLs	6
6	Internals	6
6.1	Files	6
6.2	Library	8
7	Cache	8
7.1	Structure	9
7.2	How it works	9
7.3	Administration tools	10
8	Files and directories	10
8.1	Modules	10
8.2	Configuration of the Grid Manager	11
8.3	Configuration of the GridFTP Server	16
8.4	Authorization	19
8.5	Directories	19
8.6	LRMS support	20
8.7	Runtime environment	21
9	Installation	24
9.1	Requirements	24
9.2	Setup of the Grid Manager	24
9.3	Setup of the GridFTP Server	24
9.4	Usage	25
9.5	Unix accounts	25

1 Introduction

One of the problems the user of widely distributed computing networks faces is different configuration of *Computing Elements* (CE) controlled by different administrators. This makes even initial preparation of a job non-trivial task. This is especially important in case of NorduGrid [1], where some CEs are not dedicated to NorduGrid and can not be completely reconfigured at low level. Thus some layer capable of performing most of site-dependent pre- and post-computation job is necessary.

The aim of *grid-manager* (GM) is to take care of job pre- and post-processing. It provides an interface to stage-in files containing input data and program modules from wide range of sources and transfer or store output results.

The GM is part of the NorduGrid software (codename ARC - Advanced Resource Connector). For it's connection to other parts please read "An Overview of The NorduGrid Architecture Proposal" [2]. It is **heavily using** Globus ToolkitTM 2 as it's underlying software and currently **completely depends** on it.

Essential additionally part of the GM is the specialized GridFTP Server (GFS). This server supports rich enough subset of gsiftp protocol and has network and local file access parts separated. In context of GM it's main purpose is to provide control for job and access to the job files based on the user subject and job owner. Another option is Job Control Web Service (JCS) interface implemented as part of HTTPSD framework [3].

All software described here is part of ARC software toolkit developed by NorduGrid project <http://www.nordugrid.org>

You should use this document for advanced configuration purposes. It explains internals of the aforementioned tools and extended description of configuration options. For installation and configuration refer to other documents available at <http://www.nordugrid.org/papers.html>.

2 Main concepts

A job is a set of input files (which may or may not include executables), a main executable and a set of output files. The process of gathering input files, executing a job, and transferring/storing output files is called a *session*.

Each job gets a directory on the CE called the *session directory* (SD). Input files are gathered in the SD. The job is supposed to produce new data files also in the SD. GM does not guarantee the availability of any other places accessible by the job other than SD (unless such place is part of requested Runtime Environment). The SD is also the only place which is controlled by the GM. It is accessible by the user from outside through GridFTP protocol. Any file created outside the SD is not controlled by the GM. Any exchange of data between client and GM (including also program modules) is done through GridFTP protocol [4] **only**. A URL for accessing input/output files is constructed from the base URL available through the NorduGrid Information System as part of `nordugrid-cluster` under attribute `nordugrid-cluster-contactstring` and *jobid* (jobid corresponds to a SD).

Each job gets an identifier (*jobid*). This is a handle which identifies the job in the GM and the NorduGrid Information System [5].

Each job is initiated and controlled through GFS. All job parameters (not data) are passed to the GM through GFS in RSL [6] or JSDL-coded [7] description (job description - JD). The GM adds it's own attributes to Globus RSL [8].

3 Input/output data

The main task of the GM is to take care of processing input and output data (files) of the job. Input files are gathered in SD. There are 2 ways to put file into the SD:

- Downloads initiated by the GM. Such files (name and source) are defined in the JD. It is a sole responsibility of the GM to make sure that a file will be available in the SD.
The supported sources are at the moment: GridFTP, FTP, HTTP, HTTPS (HTTP over SSLv3), HTTPg (HTTP over GSI) and SRM. Also some indexing service sources are supported: LFC, RC and RLS.
- Upload initiated by the user directly or through the User Interface (UI). Because the SD becomes available immediately at the time of submission of JD, UI can (and should) use that to upload data files which are not otherwise accessible by the GM. An example of such files can be the main executable of the job, files containing job's options/parameters, etc. These files can (and should) also be specified in the JD.

There is no other reliable way for a job to obtain input data on the CE belonging to NorduGrid. Access to AFS, NFS, FTP, HTTP and any other remote data transport during execution of the job is not guaranteed (at least not by GM).

Job stores output files in the SD. Those files also belong to 2 groups:

- Files which are supposed to be moved to a *Storage Element* (SE) and optionally registered in some *Indexing Service* like *Replica Catalog* (RC). The GM takes care of those files. They have to be specified in the JD. If job fails during any stage of processing no attempt is taken to transfer those files to their final destination, unless there is option *preserve=yes* specified in their URLs.
- Files which are supposed to be fetched by the user. The user runs UI to obtain those files. They **must** also be specified in the JD.

4 Job flow

From the point of view of the GM a job passes through various states. Picture 1 presents a diagram of the possible states of a job. A user can examine the state of a job by querying the NorduGrid Information System (IS) using the UI or any other tool. Please remember that IS can manipulate with state names to make them more user friendly and to combine them with states introduced by other parts of whole setup. Another way is to access virtual informational files through GridFTP interface or to use query method of JCS.

Configuration can put limits on amount of simultaneous jobs at some states. If such limit is reached job stays in it's current state waiting for free slot. This situation is presented by prepending current state name with **PENDING:** status mark.

Below is description of all actions taken by the GM at every state:

- **Accepted** - At this state the job has been submitted to a CE but not processed yet. The GM will analyze the JD and move to the next stage. If JD can not be processed the job will be canceled and moved to the state **Finishing**.
- **Preparing** - The input data is being gathered in the SD. The GM is downloading files specified in the JD and waiting for files which are supposed to be downloaded by the UI. If all files are successfully gathered the job moves to the next state. If **any** file can't be downloaded or it takes UI too long to upload a file - the job moves to **Finishing** state. It is possible to put limit on number of simultaneous **Preparing** jobs. Those jobs out of limit will stay in previous **Accepted** state with PENDING mark. Exceptions are jobs which has no files to be downloaded. Those are processed out of limits.

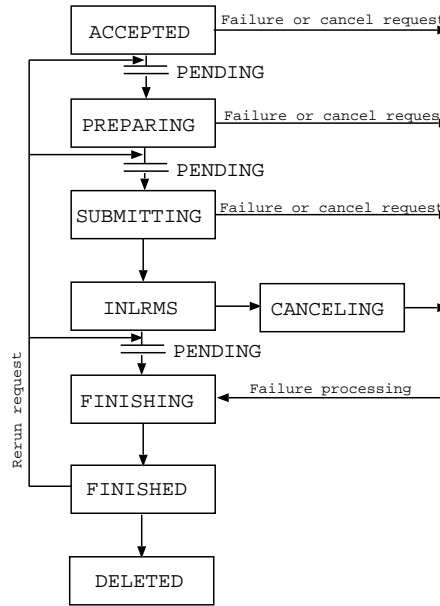


Figure 1: Job states

- **Submitting** - This is a point of interaction with *Local Resource Management System* (LRMS). At the moment PBS is supported best and corresponding backends are provided with default installation. The job is being submitted for execution. If the local job submission is successful the job moves to the next state. Otherwise it moves to **Finishing**. It is possible to limit number of jobs in **Submitting** and following **InLRMS** states.
- **InLRMS** - The job is queued or being executed in the LRMS. The GM takes no actions except waiting until job finishes.
- **Cancelling** - Necessary action to cancel job in the LRMS is being taken.
- **Finishing** - The output data is being processed. Specified data files are moved to the specified SEs and are optionally registered at RC. The user can download data files from the SD by using UI or any other tool. All the files not specified as output files are removed from the SD at very beginning of this state. It is possible to limit number of simultaneous jobs in this state.
- **Finished** - No more processing is performed by the GM. The user can continue to download data files from the SD. The SD is kept available for some time (default is 1 week). After that it is moved to the state **Deleted**. The 'deletion' time can be queried at NorduGrid Information System as attribute `nordugrid-pbs-job-sessiondireraset` of `nordugrid-pbs-job`. If job was moved to **Finished** because of failure, it may be restarted on request of client. Job is moved to a state previous to one which failed and is assigned mark **PENDING**. This is needed in order to not break the configuration limits. Exception is a job failed in **InLRMS** state and lacking input files specified in JD. Such job is treated like failed in **Preparing** state.
- **Deleted** - Job is moved to this state if user does not request job to be cleaned. Only minimal subset of information about such job is kept.

In the case of the failure special processing is applied to output files. All specified output files are treated as **downloadable by user**. No files will be moved to the SE.

5 URLs

The GM and its components support following data transfer protocols and corresponding URLs: *ftp*, *gsiftp*, *http*, *httpg*, *https*, *lfc*, *se*, *srm* (v1 and 2.2), *rc* and *rls*. For more information please see “Protocols, Uniform Resource Locators (URL) and extensions supported in ARC” document [9].

6 Internals

6.1 Files

For each local UNIX user listed in the GM configuration a *control directory* exists. In this directory the GM stores information about jobs belonging to that user. Multiple users can share the same *control directory*. To make it easier to recover in the case of failure, the GM stores most information in files rather than in memory. All files belonging to same job have names starting with **job.ID**. here ID is the job identifier.

The files in the control directory and their formats are described below:

- *job.ID.status* - current state of the job. It contains one word of text representing the current state of the job. Possible values are :
 - ACCEPTED
 - PREPARING
 - SUBMITTING
 - INLRMS
 - FINISHING
 - FINISHED
 - CANCELING
 - DELETED

See section 4 for a description of the various states. Additionally each value can be prepended with a prefix “PENDING:” (like PENDING:ACCEPTED, see section 4). That is used to show that the job is *ready* to be moved to a next state and it stays in a current state *only* because some limits set in configuration are exceeded.

- *job.ID.description* - contains the RSL description of the job.
- *job.ID.local* - information about job used by the GM. It consists of lines of format “*name = value*” . Not all of them are always available. The following names are defined:
 - *subject* - user subject also known as the distinguished name (DN)
 - *starttime* - the GMT time when the job was accepted represented in Generalized Time format of LDAP
 - *lifetime* - time period to live for the SD after job finished in seconds
 - *cleanup time* - the GMT time when job to be removed from cluster and SD deleted in Generalized Time format
 - *notify* - email addresses and flags to send mail to about job specified status changes
 - *processtime* - the GMT time when to start processing the job in Generalized Time format

- *exectime* - the GMT time when to start job execution in Generalized Time format
- *expiretime* - the GMT time when credentials delegated to job expire in Generalized Time format
- *rerun* - number of retries left to run the job
- *jobname* - name of the job as supplied by the user
- *lrms* - name of LRMS to run the job at
- *queue* - name of the queue to run the job at
- *localid* - job id in LRMS (appears only then the job is at state **InLRMS**)
- *args* - list of command-line arguments including the executable
- *downloads* - number of files to download into SD before execution
- *uploads* - number of files to upload from SD after execution
- *gmlog* - directory name which holds files containing information about job when accessed through GridFTP interface
- *clientname* - name and ip address:port of client machine (name is provided by user interface)
- *clientsoftware* - version of software used to submit job
- *sessiondir* - SD of job
- *failedstate* - state at which job failed (available only if it is possible to restart job)
- *jobreport* - URL of *logger service* used to keep track of executed jobs (one requested by user)

This file is filled partially during job submission and fully when the job moves from the **Accepted** to the **Preparing** state.

- *job.ID.input* - list of input files. Each line contains 2 values separated by a space. First value contains name of the file relative to the SD. Second value is an URL or a file description. Example:

input.dat gsiftp://grid.domain.org/dir/input_12378.dat

url - ordinary URL for gsiftp, ftp, http, https or httpg protocols with the addition of '**replica catalog url**' (RC URL) and '**replica location service url**' (RLS URL).

Each URL can contain additional options.

file description - [size][.checksum].

size - size of the file in bytes.

checksum - checksum of the file identical to the one produced by *cksum* (1).

Both size and checksum can be left out. Special kind of file description **.** is used to specify files which are **not** required to exist.

This file is used by the '**downloader**' utility. Files with '*url*' will be downloaded to the SD and files with '*file description*' will simply be checked to exist. Each time a new **valid** file appears in the SD it is removed from the list and *job.ID.input* is updated. Any external tool can thus track the process of collecting input files by checking *job.ID.input*.

- *job.ID.output* - list of output files. Each line contains 1 or 2 values separated by a space. First value is the name of the file relative to the SD. The second value, if present, is a URL. Supported URLs are the same as those supported by *job.ID.input*.

This file is used by the '**uploader**' utility. Files with *url* will be uploaded to SE and remaining files will be left in the SD. Each time a file is uploaded it is removed from the list and *job.ID.output* is updated. Files not mentioned as output files are removed from the SD at the the beginning of the **Finishing** state.

- *job.ID.failed* - the existence of this file marks the failure of the job. It can also contain one or more lines of text describing the reason of failure. Failure includes the return code different from zero of the job itself.
- *job.ID.errors* - this file contains the output produced by external utilities like **downloader**, **uploader**, script for job submission to LRMS, etc on their stderr handle. Those are not necessarily errors, but can be just useful information about actions taken during the job processing. In case of problem include content of that file while asking for help.
- *job.ID.diag* - information about resources used during execution of job and other information suitable for diagnostics and statistics. It's format is similar to that of *job.ID.local*. The following names are at least defined:
 - *nodename* - name of computing node which was used to execute job,
 - *runtimeenvironments* - used runtime environments separated by `;`,
 - *exitcode* - numerical exit code of job,
 - *frontend_distribution* - name and version of operating system distribution on frontend computer,
 - *frontend_system* - name of operating on frontend computer,
 - *frontend_subject* - subject (DN) of certificate representing frontend computer,
 - *frontend_ca* - subject (DN) of issuer of certificate representing frontend computer,and other information provided by GNU *time* utility. Note that some implementation of *time* insert unrequested information in their output. Hence some lines can have broken format.
- *job.ID.proxy* -delegated GSI proxy.
- *job.ID.proxy.tmp* - temporary GSI proxy with different unix ownership used by processes run with effective *user id* different from job owner's *id*.

There are other files with names like *job.ID.** which are created and used by different parts of the GM. Their presence in the *control directory* can not be guaranteed and can change depending on changes in the GM code.

6.2 Library

There is a library *libarcdata* distributed as part of the GM. *libarcdata* (available only if built using autotools) provides support for moving data between different URLs. It's interface can be found in Appendix B.

7 Cache

The GM can cache input files. Caching is enabled if one or more cache directories are specified in the configuration file. The GM does not cache files marked as executable in a job. Caching can also be explicitly turned off by the user for each file by using the *cache=no* option in the URL (for URL options read "Protocols, Uniform Resource Locators (URL) and extensions supported in ARC" [9]). The disk space occupied by the cache is controlled by removing old files. For more information look in section 8.2.

7.1 Structure

Cached files are stored in sub-directories under the *data* directory in each main cache directory. Filenames are constructed from an SHA-1 hash of the URL of the file, and split into subdirectories based on the two initial characters of the hash. This enables the cached files to be evenly split over a number of subdirectories. If more than one cache directory is used the initial letter of the hash also determines which cache is to be used. The algorithm is simply to use the cache directory at index *i* in the list of directories, where *i* is found from the initial letter of the hash mod the number of caches. This limits the maximum number of cache directories to 16, as SHA-1 hashes use the character set *[0-9,a-f]*. In the extremely unlikely event of a collision between two URLs having the same SHA-1 hash, caching will not be used for the second file.

There is no indexing system for the cache, and any cache filename can be easily determined from a URL by using the same hashing algorithm as the GM, e.g. the standard command line tool *sha1sum*. The cache directory used by a file can be determined from the initial letter of the hash and the order of cache directories in the configuration file. Some associated metadata (the corresponding URL and an expiry time, if available) are stored in a file with the same name as the cache file, with a *.meta* suffix.

For example, with a cache directory */cache*, the file

lfc://atlaslfc.nordugrid.org//grid/atlas/file1 is mapped to */cache/data/78/f607405ab1df6b647fac7aa97dfb6089c19fb3*

and the file */cache/data/78/f607405ab1df6b647fac7aa97dfb6089c19fb3.meta* contains the original URL and an expiry time if one is available.

At the start of a file download, the cache file is locked, so that it cannot be deleted and so that another download process cannot write the same file simultaneously. This is done by creating a file with the same name as the cache filename but with a *.lock* suffix. This file contains the process ID of the process and the hostname of the host holding the lock. If this file is present, another process cannot do anything with the cache file and must wait until the cache file is unlocked (i.e. the *.lock* file no longer exists). The lock has a timeout of one day, so that stale locks left behind by a download process exiting abnormally will eventually be cleaned up. Also, if the process corresponding to the process ID stored inside the lock is no longer running on the host specified in the lock, it is safe to assume that the lock file can be deleted.

7.2 How it works

If a job requests an input file which can be cached or is allowed to be cached, it is stored in the selected cache directory, and depending on the configuration, either the file is copied to the SD or a hard link is created in a per-job directory and a soft link is created in the SD to there. The per-job directories are in the *joblinks* subdirectory of the main cache directory. The former option is advised if the cache is on a file system which will suffer poor performance from a large number of jobs reading files on it, or the file system containing the cache is not accessible from worker nodes. The latter option is the default option. The per-job directory is only readable by the local user running the job, and the cache directory is readable only by the GM user (usually root). This means that the local user cannot access any other users' cache files. It also means that cache files can be removed without needing to know whether they are in use by a currently running job. **IMPORTANT:** If a cache is mounted from an NFS server and the GM is run by the root user, the server must have the *no_root_squash* option set for the GM host in the */etc/exports* file, otherwise the GM will not be able to create the required directories.

If the file system containing the cache is full and it is impossible to free any space, the download fails and is retried without using caching.

Before giving access to a file already in the cache, the GM contacts the initial file source to check if the user has read permission on the file. Also file creation or validity times from the original source are checked to make sure the cached file is fresh enough. If it is impossible to obtain creation and invalidation times for the file, it is invalidated 24 hours after download.

The GM checks the cache periodically. If the used space on the file system containing the cache exceeds the high water-mark given in the configuration file it tries to remove the least-recently accessed files to reduce size to the low water-mark.

7.3 Administration tools

The following tools (installed in *\$NORDUGRID_LOCATION/libexec*) exist to help with administration of the cache:

- *cache-clean* - This tool is used periodically (every 2 minutes) by the GM to keep the size of each cache within the configured limits. It removes files from the cache if the total size of the cache is greater than the configured limit. It will attempt to remove files which are not locked in order of access time, starting with the earliest, until the size is lower than the configured lower limit. If the lower limit cannot be reached (because too many files are locked, or other files outside the cache are taking up space on the file system), the tool will exit before the lower limit is reached.

cache-clean -h gives a list of options. The most useful option for administrators is *-s*, which does not delete anything, but gives summary information on the files in the cache, including information on the ages of the files in the cache.

It is not recommended to run *cache-clean* manually to clean up the cache, unless it is desired to temporarily clean up the cache with different size limits to those specified in the configuration.

- *cache-list* - This tool is used to list all files present in each cache. It simply reads through all the *.meta* files and prints to stdout a list of all URLs stored in each cache and their corresponding cache filename, one per line.

8 Files and directories

8.1 Modules

The GM consists of few separate executable modules. Those are:

- *grid-manager* - Main module. It is responsible for processing the job, moving it through states, running other modules.
- *downloader* - This is a module responsible for gathering input files in the SD. It processes the *job.ID.input* file and updates it.
- *uploader* - This module is responsible for delivering output files to the specified SEs and registration at the RC. It processes and updates the *job.ID.output* file.
- *frontend-info-collector* - Utility to gather information about frontend and to put it into *job.ID.diag* file.
- *gm-kick* - Sends signal to the GM though FIFO file to wake it up. It's used to increase responsiveness of GM.

Following modules are always run under Unix account to which user is mapped.

- *smtp-send.sh* and *smtp-send* - These are the modules responsible for sending e-mail notifications to the user. The format of the mail messages can be easily changed by editing the simple shell script *smtp-send.sh*.
- *submit-*-job* - Here * stands for the name of the LRMS. Currently supported LRMS are PBS/Torque, Condor, SGE, LoadLeveler and SLURM. Also *fork* pseudo-LRMS is supported for testing purposes. This module is responsible for the job submission to the LRMS.
- *cancel-*-job* - This one is for canceling the job, which was submitted to LRMS.
- *scan-*-job* - This shell script is responsible for notifying the GM about completion of the job. It's implementation for PBS system uses server logs to extract information about jobs. If logs are not available it uses less reliable *qstat* command for that. Other backends use different techniques.

Also there is administrator utility:

- *gm-jobs* - prints list of jobs available on cluster and amount of jobs in every state.
`gm-jobs [-h] [-l] [-u uid] [-U name]`
 -l -l - print more information about each job,
 -l-u - pretend utility is run by user with id *uid*,
 -l-l - pretend utility is run by user with name *name*.

GM comes with plugins useable for various authorization purposes (see for example description of *authplugin* command below):

- *inputcheck* - checks if all input files specified in job description are downloadable.
`inputcheck [-h] [-d debug_level] RSL_file [proxy_file]`
 -lRSL_file -file with job description,
 -lproxy_file - credentials proxy.
- *lcas* - executes LCAS plugins on credentials and returns 0 if authorization passed.
`lcas credentials description [library [db [directory]]]`
 -lcredentials - path to file with credentials to authorize,
 -ldescription - path to file with job description,
 -llibrary - path to LCAS library (full or relative to LCAS directory),
 -ldb - path to LCAS DB file (full or relative to LCAS directory),
 -ldirectory - LCAS directory.

8.2 Configuration of the Grid Manager

The GM configuration is done through single configuration file. Historically GM supports 2 kinds of configuration files. For old one it looks at following places:

- *\$NORDUGRID_LOCATION/etc/grid-manager.conf*
- */etc/grid-manager.conf*

And for new one in

- */etc/arc.conf*

The old configuration file consists of empty lines, lines containing comment (line starts from #) or configuration commands. Blank spaces in arguments must be escaped using `\` or arguments must be enclosed in `'`. Command line starts from command followed by arguments separated from command and between them by spaces.

The new configuration file can also contain empty lines and comments starting from #. It is separated into sections. Each sections starts from string containing

- *[section name/subsection name/subsubsection name]*.

Each section continues till next section or end of file. One configuration file can have commands for multiple services/modules/programs. Each service get it's own section named after it. The GM uses section *[grid-manager]*. Some services can make use of multiple subsections to reflect their internal modular structure. Commands in section *[common]* apply to all services. Command lines have format

- *name="arguments string"*.

Names are same as in old configuration file. The *argument string* consists of same arguments as in old format. And they must obey same rules.

Both files support almost same commands. Following commands are defined (examples are given for new format):

Global commands (those which affect global parameters of the GM and affect all serviced users, also described in [10]):

- *daemon=yes|no* - specifies whether the GM should go to background after started. Defaults to *yes*.
- *logfile=[path]* - specifies name of file for logging debug/informational output. Defaults to */dev/null* for daemon mode and *stderr* for foreground mode.
- *user=[uid[:gid]]* - specifies user id (and optionally group id) to which the GM must switch after reading configuration. Defaults to *not switch*.
- *pidfile=[path]* - specifies file where id of GM process will be stored. Defaults to *not write*.
- *debug=number* - specifies level of debug information. More information is printed for higher levels. Currently highest effective number is 3 and lowest 0. Defaults to 2.

All commands above are generic for every daemon-enabled server in ARC NorduGrid toolkit (like GFS and HTTPSD).

- *joblog=[path]* - specifies where to store log file containing information about started and finished jobs.
- *jobreport=[URL ... number]* - specifies that GM has to report information about jobs being processed (started, finished) to centralized service running at given *URL*. Multiple entries and multiple URLs are allowed. *number* specifies how long old records have to be kept if failed to be reported. That time is in days. Last specified value becomes effective.
- *securetransfer=yes|no* - specifies whether to use encryption while transferring data. Currently works for GridFTP only. Default is *no*. It is overridden by value specified in URL options.
- *passivetransfer=yes|no* - specifies whether GridFTP transfers are passive. Setting this option to *yes* can solve transfer problems caused by firewalls. Default is *no*.
- *localtransfer=yes|no* - specifies whether to pass file downloading/uploading task to computing node. If set to *yes* the GM won't download/upload files. Instead it composes script submitted to LRMS in way to make it do that. This requires installation of GM and Globus to be accessible from computing nodes and environment variables *GLOBUS_LOCATION* and *NORDUGRID_LOCATION* to be set accordingly. Default is *no*.

- **maxjobs**=[*max_processed_jobs* [*max_running_jobs*]] - specifies maximum number of jobs being processed by the GM at different stages:
max_processed_jobs - maximal amount of jobs being processed by GM. This does not limit amount of jobs, which can be submitted to cluster
max_running_jobs - maximal amount of jobs passed to Local Resource Management System
Missing value or -1 means no limit.
- **maxlod**=[*max_frontend_jobs* [*emergency_frontend_jobs* [*max_transferred_files*]]] - specifies maximum load caused by jobs being processed on frontend:
max_frontend_jobs - maximal amount of jobs heavily using resources of frontend (applied before moving job to PREPARING and FINISHING states)
emergency_frontend_jobs - if limit of *max_frontend_jobs* is used only by PREPARING or by FINISHING jobs aforementioned number of jobs can be moved to another state. This is used to avoid case then jobs can't finish due to big amount of recently submitted jobs.
max_transferred_files - maximal number of files being transferred in parallel by every job. Used to decrease load on not so powerful frontends.
Missing value or -1 means no limit.
- **wakeupperiod**=*time* - specifies how often for external changes are performed (like new arrived job, job finished in LRMS, etc.). *time* is a minimal time period specified in seconds. Default is 3 minutes.
- **authplugin**=*state options plugin* - specifies *plugin* (external executable) to be run every time job is going to switch to *state*. Following states are allowed: ACCEPTED, PREPARING, SUBMIT, FINISHING, FINISHED and DELETED. If exit code is not 0 job is canceled by default. *Options* consist of *name=value* pairs separated by a comma. Following *names* are supported:
timeout - specifies how long in seconds execution of the plugin allowed to last (mandatory, "*timeout*=" can be skipped for backward compatibility).
onsuccess, *onfailure* and *ontimeout* - defines action taken in each case (*onsuccess* happens if exit code is 0). Possible actions are:
pass - continue execution,
log - write information about result into logfile and continue execution,
fail - write information about result into logfile and cancel job.
- **localcred**=*timeout plugin* - specifies *plugin* (external executable or function in shared library) to be run every time job has to do something on behalf of local user. Execution of *plugin* may not last longer than *timeout* seconds. If *plugin* looks like *function@path* then function *int function(char*,char*,char*,...)* from shared library *path* is called (*timeout* is not functional in that case). If exit code is not 0 current operation will fail.
- **norootpower**=*yes/no* - if set to yes all processes involved in job management will use local identity of a user to which Grid identity is mapped in order to access filesystem at path specified in **session** command (see below). Sometimes this may involve running temporary external process.
- **allowsubmit**=[*group ...*] - list of authorization groups of users allowed to submit new jobs while "allownew=no" is active in *jobplugin.so* configuration (see below in section 8.3). Multiple commands are allowed.
- **speedcontrol**=*min_speed min_time min_average_speed max_inactivity* - specifies how long/slow data transfer is allowed to take place. Transfer is canceled if transfer rate (bytes per second) is lower than *min_speed* for at least *min_time* seconds, or if average rate is lower than *min_average_speed*, or no data is received for longer than *max_inactivity* seconds.

- **copyurl=template replacement** - specifies that URLs, starting from *template* should be accessed in a different way (most probably Unix open). The *template* part of the URL will be replaced with *replacement*. *replacement* can be either URL or local path starting from '/'. It is advisable to end template with '/
- **linkurl=template replacement [node_path]** - mostly identical to *copyurl* but file won't be copied. Instead soft-link will be created. *replacement* specifies the way to access the file from the frontend, and is used to check permissions. The *node_path* specifies how the file can be accessed from computing nodes, and will be used for soft-link creation. If *node_path* is missing - *local_path* will be used instead. Both *node_path* and *replacement* should not be URLs.

NOTE: URLs which fit into *copyurl* or *linkurl* are treated as more easily accessible than other URLs. That means if GM has to choose between few URLs from which should it download input file, these will be tried first.

Per UNIX user commands:

- **mail=e-mail_address** - specifies an email address **from** which the notification mails are sent.
- **defaultttl=ttl [ttr]** - specifies the time in seconds for the SD to be available after job finished (*ttl*) and after job was deleted (*ttr*) due to *ttl*. Defaults are 7days for *ttl* and 30 days for *ttr*.
- **defaultlrms=default_lrms_name default_queue_name** - specifies names for the LRMS and queue. Queue name can also be specified in the JD (currently it is not allowed to override used LRMS by using JD). In new configuration file this command is called **lrms**.
- **session=path** - specifies path to the directory in which the SD is created. If the path is * the default one is used - *\$HOME/.jobs* . In new configuration file this command is called **sessiondir**.
- **cachedir=path [link_path]** - specifies a directory to store cached data (see 7). Multiple cache directories may be specified by specifying multiple *cachedir* commands. Cached data will be distributed evenly over the caches. Specifying no *cachedir* command or commands with an empty path disables caching. Optional *link_path* specifies the path at which *path* is accessible on computing nodes, if it is different from the path on the GM host. If *link_path* is set to '.' files are not soft-linked, nor are per-job links created, but files are copied to the session directory. In old configuration file this command is called **cache**.
- **cachesize=high_mark [low_mark]** - specifies high and low watermarks for space used by each cache, as a percentage of the space on the file system on which the cache directory is located. When the max is exceeded, files will be deleted to bring the used space down to the min level. It is a good idea to have each cache on its own separate file system. To turn off this feature, "cachesize" without parameters can be specified. These cache settings apply to all caches specified by *cachedir* commands.
- **maxrerun=number** - specifies maximal number of times job will be allowed to rerun after it failed in LRMS. Default value is 2. This only specifies a upper limit. Actual number is provided in job description and defaults to 0.

All per-user commands should be put before *control* command which initiates serviced user.

- **control=path username [username [...]]** - This option initiates UNIX user as being serviced by the GM. *path* refers to the control directory (see section 6 for the description of control directory). If the path is * the default one is used - *\$HOME/.jobstatus* . *username* stands for UNIX name of the local user. Multiple names can be specified. If the name is * it is substituted by all names found in file */etc/grid-security/grid-mapfile* (for the

format of this file one should study the Globus project [11]).

Also the special name '.'(dot) can be used. Corresponding control directory will be used for **any** user. This option should be the last one in the configuration file. In new configuration file command **controldir=path** is also available. It uses special username '.' and is always executed last independent of placement in file.

- **helper=username command [argument [argument [...]]]** - associates external program with the local UNIX user. This program will be kept running under account of the specified user. *username* stands for the name of the user. Special names can be used: '*' - all names from /etc/grid-security/grid-mapfile, '.' - root user. The user should be already configured with *control* option (except root, who is always configured). *command* is an executable and *arguments* are passed as arguments to it.

Following are global commands supported only in new configuration file. Most of them are specific to underlying LRMS (PBS in this case) and are passed in environment variables if old configuration file is used.

- **pbs_bin_path=path** - path to directory which contains PBS commands.
- **pbs_log_path=path** - path to directory with PBS server's log files.
- **gnu_time=path** - path to *time* utility.
- **tmpdir=path** - path to directory for temporary files.
- **runtime_dir=path** - path to directory which contains *runtimeenvironment* scripts.
- **shared_filesystem=yes/no** - if computing nodes have an access to session directory through a shared filesystem like NFS. Corresponds to an environment variable `RUNTIME_NODE_SEES_FRONTEND`.
- **nodename=command** - command to obtain hostname of computing node.
- **scratchdir=path** - path on computing node where to move session directory before execution.
- **shared_scratch=path** - path on frontend where **scratchdir** can be found.
- **nodename=command** - command to obtain hostname of computing node.

In the command arguments (paths, executables, ...) following substitutions can be used:

%R	- session root - see command <i>session</i>
%C	- control dir - see command <i>control</i>
%U	- username
%u	- userid - numerical
%g	- groupid - numerical
%H	- home dir - home specified in /etc/passwd
%Q	- default queue - look command 'defaultlrms'
%L	- default lrms - look command 'defaultlrms'
%W	- installation path - \${NORDUGRID_LOCATION}
%G	- globus path - \${GLOBUS_LOCATION}

%c	- list of all control directories
%I	- job's ID (for plugins only, substituted in runtime)
%S	- job's state (for <i>authplugin</i> plugins only, substituted in runtime)
%O	- reason (for <i>localcred</i> plugins only, substituted in runtime). Possible reasons are:
new	- new job, new credentials
renew	- old job, new credentials
write	- write/delete file, create/delete directory (through gridftp)
read	- read file, directory, etc. (through gridftp)
extern	- call external program (grid-manager)

Some configuration parameters can be specified from command line while starting the GM:

grid-manager [-h] [-C level] [-d level] [-c path] [-F] [-U uid[:gid]] [-L path] [-P path]

-h - short help,

-d - debug level,

-L - name log file (overwrites value in configuration file),

-P - name for file containing process id (overwrites value in configuration file),

-U - user and group id to use for running daemon,

-F - do not make process daemon,

-c - name of configuration file,

-C - remove old information before starting: 1- remove finished jobs, 2 - remove active jobs too, 3- also remove everything that looks like junk.

8.3 Configuration of the GridFTP Server

Default location of the GFS configuration file is */etc/arc.conf* or *\$NORDUGRID_LOCATION/etc/gridftpd.conf*. Format of these configuration files is similar to that of the GM. It also supports generic commands described at the beginning of previous section 8.2. In the new format sections [common] and [gridftpd] are used. Commands specific to the GFS are described below.

- **port=number** - specifies TCP/IP port number. Default is 2811.
- **include=path** - include contents of another file. Generic commands can't be specified there.
- **encryption=yes|no** - specifies if server will allow data transfer to be encrypted. Default is yes.
- **pluginpath=path** - specifies the path where plugin libraries are installed
- **allowunknown=yes|no** - if set to *yes* clients are not checked against grid-mapfile. Hence only access rules specified in this configuration file will be applied.
- **firewall=hostname** - use IP address of the *hostname* in response to PASV command instead of IP address of a network interface of computer. You can write IP address directly instead of *hostname*. This command may be if server is situated behind NAT.

- **unixgroup=group rule** - define local UNIX user and optionally UNIX group to which user belonging to specified authorization *group* is mapped (see Section 8.4 for definition of group). Local names are obtained from specified *rule*. If specified rule could not produce any mapping, next command is used. Mapping stops at first matched rule. Following rules are supported:
 - **mapfile file** - user's subject is matched against list of subjects stored in specified file, one per line followed by local UNIX name.
 - **simplepool directory** - user is assigned one of local UNIX names stored in a file *directory/pool*, one per line. Used names are stored in other files placed in the same *directory*. If UNIX name was not used for 10 days, it may be reassigned to another user.
 - **lcm maps library directory database** - call LCMAPS functions to do mapping. Here *library* is path to shared library of LCMAPS, either absolute or relative to *directory*; *directory* is path to LCMAPS installation directory, equivalent of LCMAPS_DIR variable; *database* is path to LCMAPS database, equivalent to LCMAPS_DB_FILE variable. Each arguments except *library* is optional and may be either skipped or replaced with '*'.
 - * %D - subject of users's certificate,
 - * %P - name of credentials' proxy file.
 - **mapplugin timeout plugin [arg1 [arg2 [...]]]** - run external *plugin* executable with specified arguments. Execution of *plugin* may not last longer than *timeout* seconds. Rule matches if exit code is 0 and there is UNIX name printed on *stdout*. Name may be optionally followed by UNIX group separated by ':'. In arguments following substitutions are applied before plugin is started:
- **unixvo=vo rule** - same as **unixgroup** for users belonging to Virtual Organization (VO) *vo*.
- **unixmap=[unixname][:unixgroup] rule** - define local UNIX user and optionally group used to represent connected client. *rule* is one of those allowed for **authorization groups** (see Section 8.4) and for **unixgroup/unixvo**. In case of mapping rule username is one, provided by rule. Otherwise specified *unixname:unixgroup* is taken. Both *unixname* and *unixgroup* may be either omitted or set to '*' to specify missing value.
- **groupcfg=name** - is put into subsections representnig plugin or [group] section and defines if that section is effective. In old format it selects the group to which all following lines apply. Only unaffected option is **groupcfg**. If name is empty (or no groupcfg is used at all) following lines apply to all users.

Subsections of *gridftp* section specifies plugins which serve virtual FTP path (similar to mount command of UNIX). Name of subsection is irrelevant. In old format this section starts with command **plugin path library_name** and ends with keyword **end**. Inside subsection following commands are supported

- **plugin=library_name** - use plugin *library_name* to serve virtual path.
- **path=path** - virtual path to serve.

GFS comes with 3 plugins: *fileplugin.so*, *gacplugin.so* and *jobplugin.so*.

- *jobplugin.so* does not require any specific options in case of old configuration format. It reads the configuration file of the GM located at the standard place as specified in the section 8.2. Following options are supported:
 - * **configfile=path** - defines non-standard place for GM's configuration file,
 - * **allownew=yes|no** - specifies if new jobs can be submitted. Default is *yes*.

- * **unixgroup/unixvo/unixmap** - same options like in top-level GFS configuration. If mapping succeeds obtained local user will be used to run submitted job.

– *filepligin.so* supports following options:

- * **mount=path** - defines the place on local filesystem to which file access operations apply
- * **dir=path options** - specifies access rules for accessing files in *path* (relative to virtual and real path) and all the files below.

options is the list of the following keywords:

- **nouser** - do not use local file system rights, only use those specifies in this line
- **owner** - check only file owner access rights
- **group** - check only group access rights
- **other** - check only "others" access rights

The options above are exclusive. If none of the above specified usual Unix access rights are applied.

- **read** - allow reading files
- **delete** - allow deleting files
- **append** - allow appending files (does not allow creation)
- **overwrite** - allow overwriting already existing files (does not allow creation, file attributes are not changed)
- **dirlist** - allow obtaining list of the files
- **cd** - allow to make this directory current
- **create owner:group permissions_or:permissions_and** - allow creating new files. File will be owned by *owner* and owning group will be *group*. If '*' is used, the user/group to which connected user is mapped will be used. The permissions will be set to *permissions_or* & *permissions_and* (second number is reserved for the future usage).
- **mkdir owner:group permissions_or:permissions_and** - allow creating new directories.

- *gacpligin.so* does not have options in case of the old configuration. First line of it's configuration contains local path (root directory) served by it. Rest till keyword **end** contains GACL [12] XML used to setup initial access rules for every newly created file and directory. If GACL XML is empty then there will be no default ACLs created for new files and directories. That means ACL of parent directory will be used. For the new configuration format following options are supported: **gac=gac** - GACL XML, **mount=path** - local path server by plugin.

XML may contain variables which are replaced with values taken from client's credentials. Following variables are supported:

- \$subject** - subject of user's certificate (DN),
- \$voms** - subject of VOMS[13] server (DN),
- \$vo** - name of VO (from VOMS certificate),
- \$role** - role (from VOMS certificate),
- \$capability** - capabilities (from VOMS certificate),
- \$group** - name of group (from VOMS certificate) .

Additionally root directory must contain *.gac* file with initial ACL. Otherwise rule will be "deny all for everyone".

Some configuration parameters can be specified from command line while starting the GFS:

gridftpd [-h] [-p number] [-n number] [-b number] [-B number] [-d level] [-c path] [-F] [-U uid[:gid]] [-L path] [-P path]

- h - short help,
- d - debug level,
- L - name log file (overwrites value in configuration file),
- P - name for file containing process id (overwrites value in configuration file),
- U - user and group id to use for running daemon,
- F - do not make process daemon,
- c - name of configuration file,
- p - TCP/IP port number,
- n - maximal number of simultaneously served connection,
- b - default size of buffer used for data transfer (default is 64kB),
- B - maximal size of buffer used for data transfer (default is 640kB).

8.4 Authorization

Authorization is performed at GFS by applying set of rules. Each rule takes one line in the *group* section. For information about supported rules please read “Configuration and authorisation of ARC (Nordugrid) Services” [10].

8.5 Directories

The GM is installed into a single installation point referred as `$NORDUGRID_LOCATION` and following sub-directories are used:

`$NORDUGRID_LOCATION/bin` - tools

`$NORDUGRID_LOCATION/libexec` - program modules used by GM

`$NORDUGRID_LOCATION/etc` - configuration files, deprecated, central configuration file is used by default

`$NORDUGRID_LOCATION/sbin` - daemons

`$NORDUGRID_LOCATION/lib` - gridftp server's plugins and API libraries

The GM also uses following directories:

- *session root directory* - In this directory the SD is created. It can be multiple directories for the various users specified in the configuration file.

There are 2 processes which need to have a permissions to create new files and directories in it. Those are GM and GFS.

If any of those processes are run under dedicated user account, that account need full permissions in the *session root directory*.

If those processes are run under *root* account make sure *session root directory* is not on filesystem which limits capabilities of *root* user. For example NFS with *root_squash* option.

If there is need to run processes under *root* account (to run jobs in LRMS under different users' accounts) but there is no way to provide suitable *session root directory* use *norootpower* command in configuration of the GM. In that case GM and GFS will use identity of local user to which Grid identity is mapped to access *session root directory*. Hence those users will need full access there.

The GM creates SD with proper ownership and permissions for local identity used to run job. Some filesystems require *executable* permissions on *session root directory* to be set for local identity in order to access any file or subdirectory there.

This directory should also be shared among cluster nodes in order for job to access input files. Or internal means of LRMS must be used to transfer files to executing node. For more see section 8.6.

- *control directory* - In this directory the SD stores an information about the accepted jobs. Both GM and GFS processes must have full permissions there.
Also subdirectory called *log* is created there. It is used to accumulate information about started and finished jobs. This information is periodically sent to the *logger service*.

8.6 LRMS support

The GM comes with support for several LRMS. And this number is slowly growing. Features explained below are for **PBS** backend. This support is provided through *submit-pbs-job*, *cancel-pbs-job*, *scan-pbs-job* scripts. *submit-pbs-job* creates job's script and submits it to PBS. Created job's script is responsible for moving data between frontend machine and cluster node (if required) and execution of actual job. Alternatively it can download input files and upload output if "*localtransfer no*" is specified in the configuration file.

Behavior of submission script is mostly controlled using environment variables. Most of them can be specified on frontend in GM's environment and overwritten on cluster's node through PBS configuration. Some of them may be set in configuration file too.

PBS_BIN_PATH - path to PBS executables. Like */usr/local/bin* for example. *pbs_bin_path* configuration command.

PBS_LOG_PATH - path to PBS server logs. *pbs_log_path* configuration command.

TMP_DIR - path to directory to store temporary files. Default value is */tmp*. *tmpdir* configuration command.

RUNTIME_CONFIG_DIR - path where runtime setup scripts can be found. *runtimedir* configuration command.

GNU_TIME - path to GNU time utility. It is important to path to utility compatible with GNU time. If such utility is not available, modify *submit-pbs-job* to either reset this variable or change usage of available utility. *gnu_time* configuration command.

NODENAME - command to obtain name of cluster's node. Default is */bin/hostname -f*. *nodename* configuration command.

RUNTIME_LOCAL_SCRATCH_DIR - if defined should contain path to the directory on computing node, which can be used to store job's files during execution. *scratchdir* configuration command.

RUNTIME_FRONTEND_SEES_NODE - if defined should contain path corresponding to **RUNTIME_LOCAL_SCRATCH_DIR** as seen on **frontend** machine. *shared_scratch* configuration command.

RUNTIME_NODE_SEES_FRONTEND - if set to "*no*" means computing node does not share filesystem with frontend. In that case content of the SD is moved to computing node by using means provided by the LRMS. Results are moved back after job's execution in a same way. *shared_filesystem* configuration command.

Figures 2,3,4 present some possible combinations for **RUNTIME_LOCAL_SCRATCH_DIR** and **RUNTIME_FRONTEND_SEES_NODE** and explain how data movement is performed. Pictures a) correspond to situation right after all input files are gathered in session directory and actions taken right after job's script starts. Pictures b) show how it looks while job is running and actions which are taken right after it finished. Pictures c) stand for final situation, when job files are ready to be uploaded to external storage element or be downloaded by user.

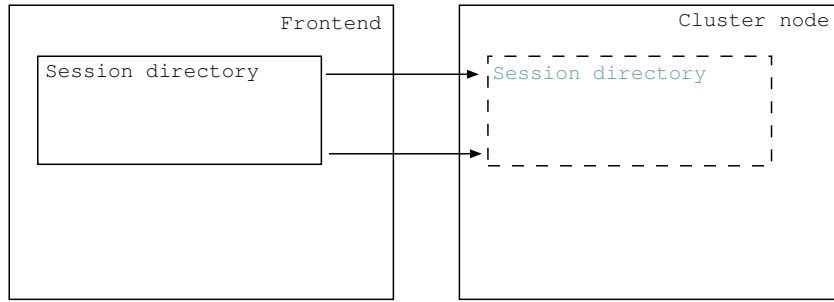


Figure 2: Both `RUNTIME_LOCAL_SCRATCH_DIR` and `RUNTIME_FRONTEND_SEES_NODE` undefined. Job is executed in session directory placed on frontend.

8.7 Runtime environment

The GM can run specially prepared *BASH* scripts prior creation of job's script, before and after executing job's main executable. Those scripts are requested by user through *runtimeenvironment* attribute in RSL and are run with only argument set equal to '0', '1' or '2' during creation of job's script, before execution of main executable and after main executable finished accordingly. They all are run through BASH's 'source' command, and hence can manipulate with shell variables. With argument '0' scripts are run by the GM on frontend. Some environment variables are defined in that case and can be changed to influence job's execution later:

- `joboption_directory` - session directory.
- `joboption_arg_#` - command and arguments to be executed as specified in RSL.
- `joboption_env_#` - array of 'NAME=VALUE' environment variables (**not** bash array).
- `joboption_runtime_#` - array of requested *runtimeenvironment* names (**not** bash array).
- `joboption_num` - *runtimeenvironment* currently beeing processed (number starting from 0).
- `joboption_stdin` - name of file to be attached to stdin handle.
- `joboption_stdout` - same for stdout.
- `joboption_stderr` - same for stderr.
- `joboption_cputime` - amout of CPU time requested (minutes).
- `joboption_memory` - amout of memory requested (megabytes).
- `joboption_count` - number of processors requested.
- `joboption_lrms` - LRMS to be used to run job.
- `joboption_queue` - name of a queue of LRMS to put job into.
- `joboption_nodeproperty_#` - array of properties of computing nodes (LRMS specific, **not** bash array).
- `joboption_jobname` - name of the job as given by user.
- `joboption_rsl` - whole RSL for very clever submission scripts.
- `joboption_rsl_name` - RSL attributes and values (like `joboption_rsl_executable="/bin/echo"`)

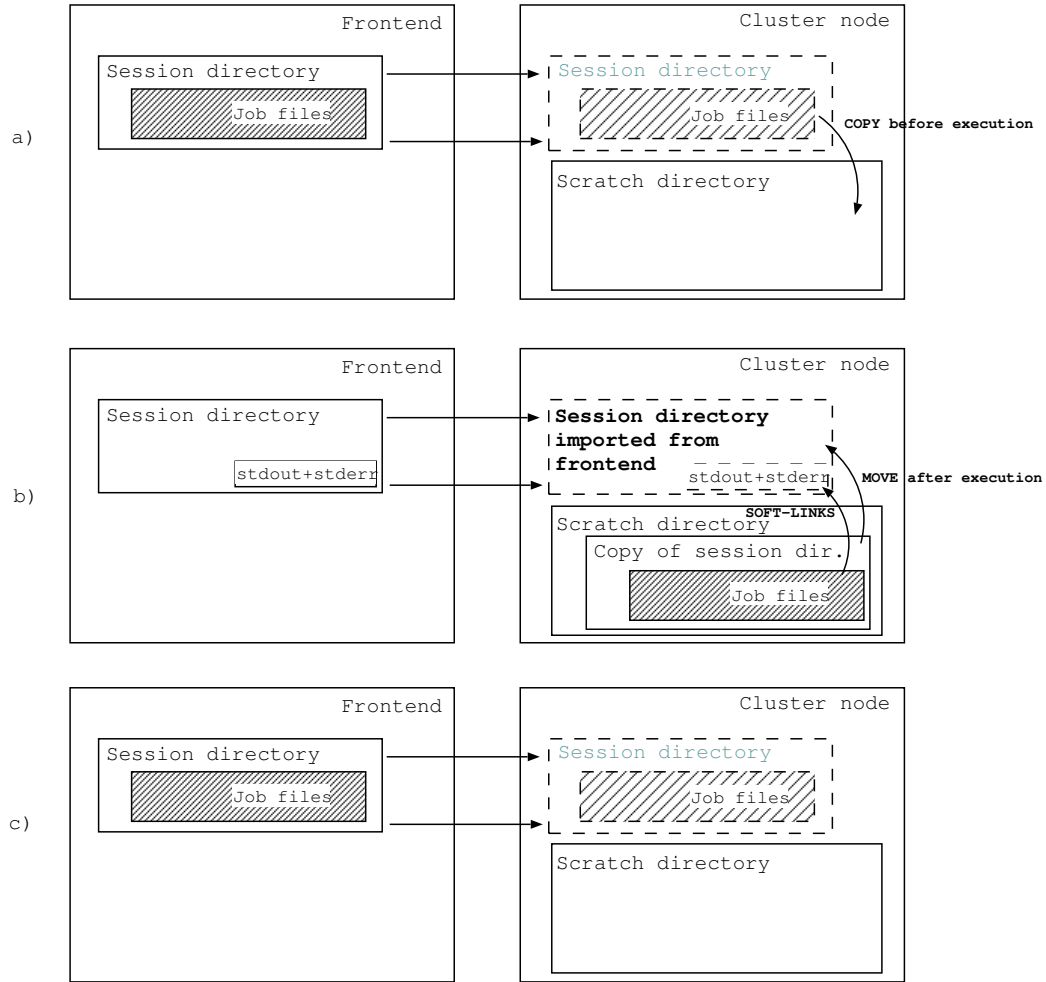


Figure 3: `RUNTIME_LOCAL_SCRATCH_DIR` is set to value representing scratch directory on computing node, `RUNTIME_FRONTEND_SEES_NODE` undefined.

- a) After job script starts all input files are moved to 'scratch directory' on computing node.
- b) Job runs in separate directory in 'scratch directory'. Only files representing job's *stdout* and *stderr* are placed in original 'session directory' and soft-linked in 'scratch'. After execution all files from 'scratch' are moved back to original 'session directory'.
- c) All output files are in 'session directory' and are ready to be uploaded/downloaded.

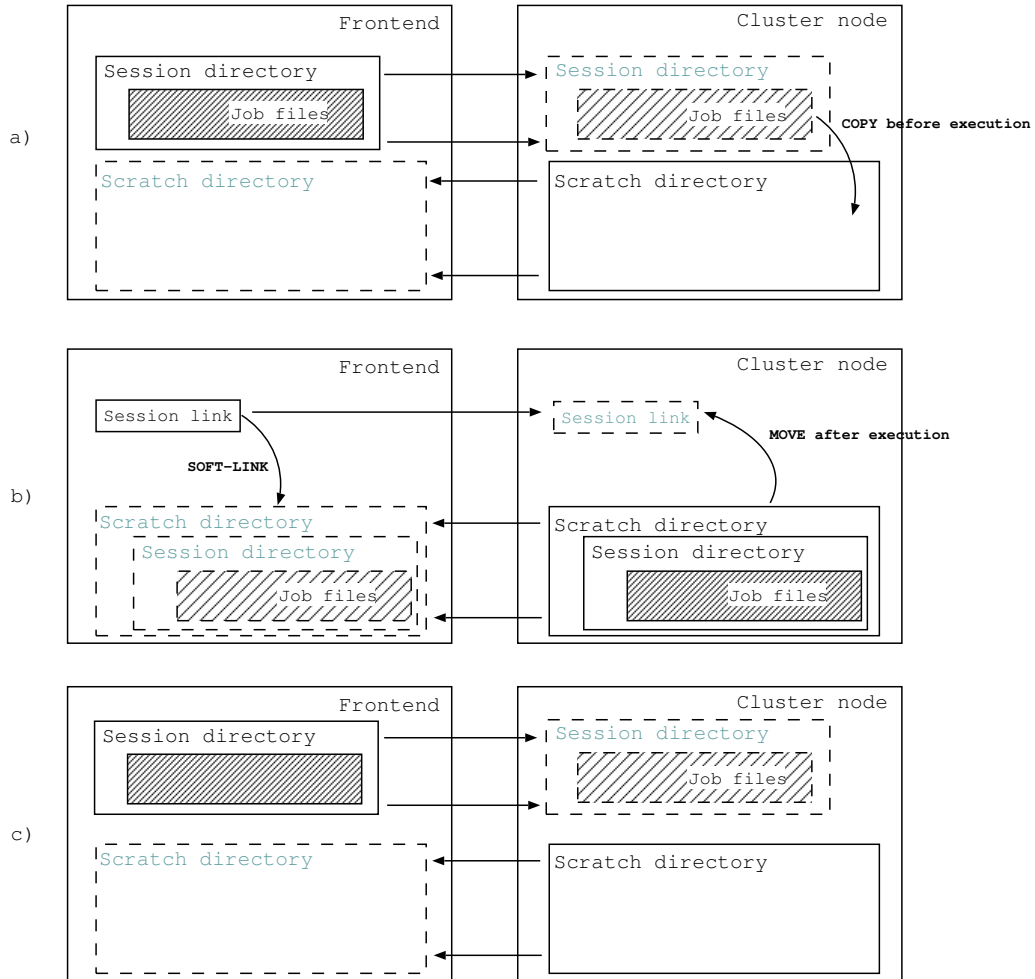


Figure 4: Both `RUNTIME_LOCAL_SCRATCH_DIR` and `RUNTIME_FRONTEND_SEES_NODE` are set to values representing scratch directory on computing node and way to access that scratch from frontend correspondingly.

- After job script starts all input files are moved to 'scratch directory' on computing node. Original 'session directory' is removed and replaced with soft-link to copy of session directory in 'scratch' as seen on frontend.
- Job runs in separate directory in 'scratch directory'. All files are also available on frontend through soft-link. After execution soft-link is replaced with directory and all files from 'scratch' are moved back to original 'session directory'.
- All output files are in 'session directory' and are ready to be uploaded/downloaded.

For example *joboption_args* could be changed to wrap main executable. Or *joboption_runtime* could be expanded if current one depends on others.

With argument '1' scripts are run just before main executable is run. They are executed on computing node. Such script can prepare environment for some third-party software package. A current directory in that case is one which would be used for execution of job. Variable HOME also points to that directory.

With argument '2' scripts are executed after main executable finished. Main purpose is to clean possible changes done by scripts run with '1' (like removing temporary files). Execution of scripts at that stage also happens on computing node and is not reliable. If the job is killed by LRMS they most probably won't be executed.

9 Installation

To install GM as part of ARC-enabled site please read “Nordugrid ARC server installation instructions” at <http://www.nordugrid.org/documents/ng-server-install.html>.

9.1 Requirements

The GM is mostly written using C++. It was tested and should compile on recent enough *Linux* systems using *gcc* compiler and *GNU make* (gcc versions 2.95, 2.96, 3.2, 3.4 were tested). You will also need *Globus ToolkitTM* of version higher than 2.2 installed <http://www-unix.globus.org/toolkit/>.

9.2 Setup of the Grid Manager

For in-depth information about how to properly setup the GM and related software please read “Nordugrid ARC server installation instructions” at <http://www.nordugrid.org/documents/ng-server-install.html>. Follow that manual to install GM, configure and run it. Additional tips are described here.

The GM is designed to be able to run both as root and as ordinary user. You can chose the name of the user by using corresponding command in configuration file. It is better run GM as root if You want to serve few users.

The GM writes debug information into a file */var/log/grid-manager.log* by default. . Also file */var/log/gm-jobs.log* (default path in configuration template, turned off by default) contains information about all started and finished jobs, 2 lines per job (1 when job is started and 1 after it finished).

9.3 Setup of the GridFTP Server

For in-depth information about how to properly setup the GM and related software please read “Nordugrid ARC server installation instructions” at <http://www.nordugrid.org/documents/ng-server-install.html>. Follow that manual to install GM, configure and run it. Additional tips are described here.

Local file access in the GFS is implemented through plugins (shared libraries). There are 3 plugins provided with the GFS: *fileplugin.so*, *gacplugin.so* and *jobplugin.so* . The *fileplugin.so* is intended to be uses for plain file access with the configuration senitive to user subject and is not necessary for setting a Nordugrid compatible site. The *gacplugin.so* uses GACL (<http://www.gridpp.ac.uk/authz/gacl/>) to control access to local file system. The *jobplugin.so* is using information about jobs being controlled by GM and provides access to session directories of the jobs owned by user. It also provides an interface (virtual directory and virtual operations) to submit, cancel, clean, renew credentials and obtain information about the job.

To make GFS to interoperate with other parts of the ARC only one *jobplugin.so* needs to be configured. It is advisable to use the template configuration file. You can leave only part which configures *jobplugin.so* plugin.

9.4 Usage

Refer to the description of the *User Interface* part [14] and extensions to RSL [8] for using the GM.

9.5 Unix accounts

Both GM and GFS are designed to be run by *root* UNIX account and serve multiple local UNIX and global Grid identities. Nevertheless it is possible to use *non-root* accounts to run those services. Although this means some functionality loss described below.

There are no implications on running GFS with *gacplugin* or *fileplugin* under *non-root* account as long as only Grid identity of user is used and all served files and directories are owned by server's account.

For combination of GM and GFS with *jobplugin* both services must be run either by same account or one of services must be run under *root* account. That is needed because services communicate over local filesystem, hence must have *full* access to same set of files.

As long as GFS with *jobplugin* is run under non-root account there is no mapping from Grid identity to local UNIX account taking place. All allowed Grid users are assigned server's account and are then processed by GM using same account. Only way to overcome this limitation is to run one GFS per local account with proper access control configured.

Because GM has to represent user's local account while communication with LRMS, it can serve only account it is run under (unless it is run under *root* account, of course). Like in case of GFS, multiple instances of GM may be run, one per local account. That solution causes another implications. The GM loses possibility to share cached files among serviced users. It is not also possible to control load on a frontend by limiting number of simultaneously running *downloader* and *uploader* modules.

One has also take into account that private part of GSI infrastructure (private key of a host at least) has to be duplicated for every account used to run GFS.

Appendix A. Job control over jobplugin.so

Virtual tree

Under mount point of jobplugin gridftp client can see directories representing job belonging to user, who started client. Directory per job. Directory names are same as jobs' identifiers. Those directories are directly connected to session directories of jobs and contain same files and subdirectories. Except if jobs session directory is moved to computing node. In that case directories only contain files with redirected stdout and stderr as specified in xRSL.

If job's xRSL has *gmlog* specified job's directory also contains virtual subdirectory with same name, which contains files with information about job as created by GM. The most important are 'errors' and 'status'. 'errors' contains stderr of separate modules run by GM in order to process job (downloader, uploader, job's submission to LRMS). 'status' contains one word representing state of job.

Also under mount point there is additional directory named "new" used to submit new jobs. And another directory "info" with subdirectories named after job ids. Those subdirectories contain files with information about job identical

to those in subdirectory specified through *gmlog*.

Submission

Each xRSL put into directory "new" is accepted as job's description. *jobplugin* parses it and client gets positive response if there are no errors in request.

Job gets identifier and directory with corresponding name appears. If job's description contains input files which should be delivered from client's machine, client must upload them to that directory under specified names.

Because each job gets identifier there should be a way for client to obtain it. For that prior to providing xRSL client sends command CWD to change current directory to "new". In this way job's identifier is reserved, new directory corresponding to that identifier is created and client is redirected to it (as specified in FTP protocol). Job's description put into "new" will get reserved identifier.

Actions

Various actions to affect processing of existing job are performed by uploading xRSL files into directory "new". Content of xRSL may consist of only 2 parameters - action for *action* to be performed, and *jobid* to identify job to be affected. Rest of parameters are ignored.

Currently supported actions are:

<i>cancel</i>	to cancel job
<i>clean</i>	to remove job from computing resource
<i>renew</i>	to renew credentials delegated to job
<i>restart</i>	to restart job after failure at some phases

It is also possible to perform some actions by using shortcut FTP operations described below.

Cancel

Job is canceled by performing DELE (delete file) command on directory representing job. It can take some time (few minutes) before job is actually canceled. Nevertheless client gets response immediately.

Clean

Job's content is cleaned by performing RMD (remove directory) command on directory representing job. If job is in "FINISHED" state it will be cleaned immediately. Otherwise it will be cleaned after it reaches state "FINISHED".

Renew

If client requests CWD to session directory credentials passed during authentication are compared to current credentials of the job. If validity time of the new credentials is longer job's credentials are replaced with new.

Appendix B. Library *libarcdata*

libarcdata is now part of *libngui* library. It's functions are declared in a header file *arcdata.h*. They correspond to *ng** utilities meant for data handling - *arcacl*, *arccp*, *arcls*, *arcrm*, *arctransfer*. It consists of following functions:

```
void arcacl(const std::string& file_url, const std::string& command, int timeout = 0);
void arcregister (const std::string& source_url, const std::string& destination_url, int timeout = 0);
void arccp (const std::string& source_url, const std::string& destination_url, bool recursive = false, int timeout = 0);
void arcls(const std::string& dir_url, bool show_details = false, bool show_urls = false, int timeout = 0);
void arcrm(const std::string& file_url, bool errcont = false, int timeout = 0);
void arctransfer(const std::string& destination, std::list<std::string>& sources, int timeout = 0);
```

Additionally this library contains C++ classes used by *ng** data management utilities. Those are described in “ARC::DataMove Reference Manual”.

Appendix C. Error messages of GM

If job has not finished successfully the GM put one or more lines into *job.ID.failed*. Possible values include those generated by the GM itself:

<i>Error string</i>	<i>Reason/description</i>
Internal error	Error in internal algorithm
Internal error: can't read local file	Error manipulating files in the control directory
Failed reading local job information	-//-
Failed reading status of the job	-//-
Failed writing job status	-//-
Failed during processing failure	-//-
Serious troubles (problems during processing problems)	-//-
Failed initiating job submission to LRMS	Could not run backend executable to pass job to LRMS
Job submission to LRMS failed	Backend executable supposed to pass job to LRMS returned non-zero exit code
Failed extracting LRMS ID due to some internal error	Output of Backend executable supposed to contain local ID of passed job could not be parsed
Failed in files upload (post-processing)	Failed to upload some or all output files
Failed in files upload due to expired credentials - try to renew	Failed to upload some or all output files most probably due to expired credentials (proxy certificate)
Failed to run uploader (post-processing)	Could not run <i>uploader</i> executable
uploader failed (postprocessing)	Generic error related to <i>uploader</i> component
Failed in files download (pre-processing)	Failed to upload some or all input files
Failed in files download due to expired credentials - try to renew	Failed to download some or all input files most probably due to expired credentials (proxy certificate)
Failed to run downloader (pre-processing)	Could not run <i>downloader</i> executable
downloader failed (preprocessing)	Generic error related to <i>downloader</i> component
User requested to cancel the job	GM detected external request to cancel this job, most probably issued by user
Could not process RSL	Job description could not be processed to syntax errors or missing elements
User requested dryrun. Job skipped.	Job description contains request not to process this job
LRMS error: (CODE) DESCRIPTION	LRMS returned error. CODE is replaced with numeric code of LRMS, and DESCRIPTION with textual description
Plugin at state STATE failed: OUTPUT	External plugin specified in GM's configuration returned non-zero exit code. STATE is replaced by name of state to which job was going to be passed, OUTPUT by textual output generated by plugin.
Failed running plugin at state STATE	External plugin specified in GM's configuration could not be executed.

Provided by downloader component (URL is replaced by source of input file, FILE by name of file):

<i>Error string</i>	<i>Reason/description</i>
Internal error in downloader	Generic error
Input file: URL - unknown error	Generic error
Input file: URL - unexpected error	Generic error
Input file: URL - bad source URL	Source URL is either malformed or not supported
Input file: URL - bad destination URL	Shouldn't happen
Input file: URL - failed to resolve source locations	File either not registred or other problems related to Data Indexing service.
Input file: URL - failed to resolve destination locations	Shouldn't happen
Input file: URL - failed to register new destination file	Shouldn't happen
Input file: URL - can't start reading from source	Problems related to accessing instance of file at Data Storing service.
Input file: URL - can't read from source	-//-
Input file: URL - can't start writing to destination	Access problems in a session directory
Input file: URL - can't write to destination	-//-
Input file: URL - data transfer was too slow	Timeouted while trying to download file
Input file: URL - failed while closing connection to source	Shouldn't happen
Input file: URL - failed while closing connection to destination	Shouldn't happen
Input file: URL - failed to register new location	Shouldn't happen
Input file: URL - can't use local cache	Problems with GM cache
Input file: URL - system error	Operating System returned error code where unexpected
Input file: URL - delegated credentials expired	Access to source requires credentials and they are either outdated or missing (not delegated).
User file: FILENAME - Bad information about file: checksum can't be parsed.	In job description there is a checksum provided for file uploadable by user interface and this record can't be interpreted.
User file: FILENAME - Bad information about file: size can't be parsed.	In job description there is a size provided for file uploadable by user interface and this record can't be interpreted.
User file: FILENAME - Expected file. Directory found.	Instead of file uploadable by user interface GM found directory with same name in a session directory.
User file: FILENAME - Expected ordinary file. Special object found.	Instead of file uploadable by user interface GM found special object with same name in a session directory.
User file: FILENAME - Delivered file is bigger than specified.	The size of file uploadable by user interface is bigger than specified in job description.
User file: FILENAME - Delivered file is unreadable.	GM can't check user uploadable file due to some internal error. Most probably ²⁹ due to improperly configured local permissions.
User file: FILENAME - Could not read file to compute checksum.	GM can't read user uploadable file due to some internal error. Most probably due to improperly configured local permissions.

Provided by uploader component (URL is replaced by destination of output file) :

<i>Error string</i>	<i>Reason/description</i>
Internal error in uploader	Generic error
Output file: URL - unknown error	Generic error
Output file: URL - unexpected error	Generic error
User requested to store output locally URL	Destination is URL of type <i>file</i> .
Output file: URL - bad source URL	Shouldn't happen
Output file: URL - bad destination URL	Destination URL is either malformed or not supported
Output file: URL - failed to resolve source locations	Shouldn't happen
Output file: URL - failed to resolve destination locations	Problems related to Data Indexing service.
Output file: URL - failed to register new destination file	-//-
Output file: URL - can't start reading from source	User request to store output file, but there is no such file or there are problems accessing session directory
Output file: URL - can't start writing to destination	Problems with Data Storing services
Output file: URL - can't read from source	Problems accessing session directory
Output file: URL - can't write to destination	Problems with Data Storing services
Output file: URL - data transfer was too slow	Timeout during transfer
Output file: URL - failed while closing connection to source	Shouldn't happen
Output file: URL - failed while closing connection to destination	Shouldn't happen
Output file: URL - failed to register new location	Problems related to Data Indexing service.
Output file: URL - can't use local cache	Shouldn't happen
Output file: URL - system error	Operating System returned error code where unexpected
Output file: URL - delegated credentials expired	Access to destination requires credentials and they are either outdated or missing (not delegated).

Coming from LRMS (PBS) backend:

<i>Error string</i>	<i>Reason/description</i>
Submission: Configuration error.	
Submission: System error.	
Submission: Job description error.	
Submission: Local submission client behaved unexpectedly.	
Submission: Local submission client failed.	

References

- [1] "The NorduGrid Collaboration," Web site. [Online]. Available: <http://www.nordugrid.org>
- [2] A. Wäänänen, "An Overview of an Architecture Proposal for a High Energy Physics Grid," in *Proc. of PARA 2002, LNCS 2367*, p. 76, J. Fagerholm, Ed. Springer-Verlag Berlin Heidelberg, 2002.
- [3] A. Konstantinov, *The HTTP(s,g) And SOAP Framework*, The NorduGrid Collaboration, NORDUGRID-TECH-9. [Online]. Available: http://www.nordugrid.org/documents/HTTP_SOAP.pdf
- [4] W. Allcock et al., "Data management and transfer in high-performance computational grid environments," *Parallel Comput.*, vol. 28, no. 5, pp. 749-771, 2002.
- [5] B. Kónya, *The NorduGrid/ARC Information System*, The NorduGrid Collaboration, NORDUGRID-TECH-4. [Online]. Available: http://www.nordugrid.org/documents/arc_infosys.pdf
- [6] "The Globus Resource Specification Language RSL v1.0." [Online]. Available: http://www-fp.globus.org/gram/rsl_spec1.html
- [7] A. Anjomshoaa et al., "Job Submission Description Language (JSDL) Specification v1.0," November 2007, GFD-R.056. [Online]. Available: <http://www.gridforum.org/documents/GFD.56.pdf>
- [8] O. Smirnova, *Extended Resource Specification Language*, The NorduGrid Collaboration, NORDUGRID-MANUAL-4. [Online]. Available: <http://www.nordugrid.org/documents/xrsl.pdf>
- [9] A. Konstantinov, *Protocols, Uniform Resource Locators (URL) and Extensions Supported in ARC*, The NorduGrid Collaboration, NORDUGRID-TECH-7. [Online]. Available: <http://www.nordugrid.org/documents/URLs.pdf>
- [10] ---, *Configuration and Authorisation of ARC (NorduGrid) Services*, The NorduGrid Collaboration, NORDUGRID-TECH-6. [Online]. Available: http://www.nordugrid.org/documents/Config_Auth.pdf
- [11] I. Foster and C. Kesselman, "Globus: A Metacomputing Infrastructure Toolkit," *International Journal of Supercomputer*

- Applications*, vol. 11, no. 2, pp. 115-128, 1997, available at:
<http://www.globus.org>.
- [12] A. McNab, "The GridSite Web/Grid security system: Research Articles," *Softw. Pract. Exper.*, vol. 35, no. 9, pp. 827-834, 2005.
- [13] R. Alfieri *et al.*, "From gridmap-file to VOMS: managing authorization in a Grid environment," *Future Gener. Comput. Syst.*, vol. 21, no. 4, pp. 549-558, 2005.
- [14] M. Ellert, *The NorduGrid toolkit user interface*, The NorduGrid Collaboration, NORDUGRID-MANUAL-1. [Online]. Available:
<http://www.nordugrid.org/documents/ui.pdf>