

DRMAA Interface Specification

by
DRMAA Working Group participants

Document maintained by Hrabri.Rajic@intel.com

History:

Date	Document	Comment
Apr 29, 2002	Sched-drmaa-1.2	This version combines sched-drmaa-1.0 and 1.1 documents and DRMAA working group discussion till April 16.
May 8, 2002	Sched-drmaa-1.3	Addition of “job categories” and “explicit native resource specification” sections contributed by Andreas Haas. Minor polishes in response to DRMAA telecom discussion of Apr 30, ‘02.

1 Introduction

Distributed Resource Management Application API (DRMAA) hides the differences of the Distributed Resource Management Systems (DRMSs) and provides an API intended for distributed application developers or ISVs. It is one of the DRMAA working group goals to target a very broad audience and consequently require an easy learning curve for the use of the specified API. The mandate of the DRMAA working group is to produce DRMAA specification 1.0.

Language issues

In is our position that the API should be implemented in multiple languages, C/C++ being the primary choice. The secondary choices are scripting languages, Perl and Python. Perl is especially heavily used in the biotech arena where there is great need for numerous parametric calculations.

It is possible to design an Interface Definition Language that will effectively resolve the issue of one interface serving multiple languages. While this is a viable approach, we feel that it will slow the progress on the implementation side significantly. Another viable approach would be to design a protocol instead of the API. Both of these alternatives should be considered in more detail when the time comes to address the long term comprehensive solutions.

Library Issues

An ideal library would have paths to handle all DRMSs and versions to be linked statically or dynamically. This is not something that will be feasible. A real possibility is a situation where one vendor implements multiple, but not all DRMSs. The packaging could come as one library, where a DRMS is selected at run time setting an environmental variable for the desired DRMS, or as one DRMS link per library. The latter approach is advocated by the authors. In this setup the shared library is selected at run time by end users.

It is expected that the developers will be linking the library from serial and multithreaded codes. The library should be thread safe.

It is expected that the debugging of the distributed programs will be more challenging than single machine versions. We advocate providing production and debugging version of the library.

Library should provide the DRMAA API version number to the external programs (such as SCCS’s “what” and RCS’s “ident”) and to the distributed applications programmatically.

User program and DRMAA interaction

All of the DRMSs are asynchronous in nature. They notify the end user of the status of a finished job via e-mail, which as an only option is not acceptable to the users of DRMAA library. We propose to deal with the asynchrony similarly to Unix and Windows process interfaces, by blocking on the wait call for a specific job request or possibly to all of them in the same process. This is in contrast to Globus GRAM interface that is modeled on the reactive mode of execution. The support for reactive mode is to be addressed in the future DRMAA versions, because more and more programs come with graphic user interface these days.

The rest of the document is presented as follows. The Chapter 2 deals with the API design issues. Chapter 3. presents a Draft API specification.

2 API Design Issues

Developers have been using Unix system, fork/exec, popen, and the wait interfaces for years to spawn additional processes and wait for the end of their execution to get their exit codes. Windows has equivalent utilities like CreateProcess and WaitForSingleObject. DRMAA provides its own set of interfaces that are OS neutral. It borrows Unix process API simplicity and tries to be consistent with libc interfaces.

2.1 Basic Guidelines

Even though the API should be self-contained, it is not always possible to consolidate all variations of end user and DRMS interactions under the API. For this reason, we advocate that the developers should provide a way for the end user to specify DRMS particular options. The primary mean to achieve this is through use of “job categories”.

Additionally, there might be a need for possibility that the DRMS specific options are specified at run time as command line options. The end user could loose portability this way, which is a small price to pay to be able to run the application in uncommon configurations. Besides that, DRMAA providers and ISVs are the ones that target multiple DRMSs; the end user does that at a much lesser degree.

The API centers around `job_id` parameter that is passed back by the DRMS upon job submission. `Job_id` is used for all the job control and monitoring purposes. [An additional similar parameter, `job_name`, that is found in all DRMS implementations is part of the job submission interface. `Job_name` could be used by the developer and/or internally by the implementation to group the jobs for easier user classification and tracking. This parameter could be a key to achieve scalability for DRMAA implementations, especially since DRMS user jobs could be running concurrently with those of the other DRMS users.]

There are few guidelines that were used in designing the uniform API:

- The API calling sequences should be simple and the API set small.
- The routine names should convey the semantic of the routine.
- The set should be as convenient as possible, even with the risk of being forced to emulate some functionality if missing from a DRMS.
- All job manipulation is available without explicit job iterating.
- The server names are hidden, the DRMS is a black box.
- The end user does not need to interact with the DRMS since the DRMS environment is specified with job categories by the application administrator. If he has to he could specify native resource options parameter.
- The API should be extensible.

2.2 DRMAA Distributed Application Environment

DRMAA specification 1.0 does not have explicit file staging mechanisms. File staging is enabled by setting job template attributes.

2.2.1 Job categories

The DRMAA interface specification should allow ISVs to write DRM-enabled applications even though the properties of a concrete DRM installation, in particular the configuration of the DRM system, cannot be known in advance.

Experiences made with integrations based on DRM CLI show that even when the same ISV application is run as a job with the same DRM system the site specific policies in effect differ widely. These policies are typically about questions like

- what resources are to be used by the job
- preferences where to run the job
- how prior the job should be treated by the DRM scheduler compared to other jobs

For supporting the variety of policies, job specific requests expressed by DRM submit options are very common in the DRM product space.

Despite of these differences between two sites with the "same" job passed to the DRM system the application actually does not change when seen from the perspective of the ISV. Also for the end user who just wants a job to be started nothing changes due to different policies. This is an indication that there must be possibility for hiding these site-specific differences behind the DRMAA interface.

The job "categories concept" is the approach the DRMAA working group recommends for encapsulating site-specific details and completely hiding these details from applications making use of the DRMAA interface. The core of the idea is to have these application only supplying a string attribute specifying a job category, i.e. a name specifying what kind of application that is to be dispatched by the DRMS. The category name can be used by the DRMAA library to determine site specific resource and functional requirements of jobs in this category. Such requirements need to be configurable by the site operating a DRM system and deploying an ISV application on top of it.

An example can help to illustrate this idea:

- At site A rendering application X is used in a heterogeneous clustered environment which is managed by a DRMS. Since application X is only available at a subset of these machines the administrator sets up the DRMS in a way requiring from the end-users to put a **-l X=true** into their submit command line.
- At site B the same application is used in a homogenous clustered environment with rendering application X supported at all machines managed by the DRMS. However since X jobs do compete with applications Y sharing the same resources and X applications are to be treated with higher priority than Y jobs end-users need to put a **-p 1023** into their submit command line for raising the dispatch priority.

An integration based on categories will allow to submit X jobs through the DRMAA interface in compliance with the policies of both sites A and B without the need to know about these policies. The ISV does this by specifying "X" as the category used for X rendering jobs submitted through the DRMAA interface and by mentioning this in the "DRM integration" section of the X rendering software documentation.

The administrators at the sites A and B site read the documentation or installation instructions about the "X" DRMAA category. The documentation of their DRMS contains directions about the category support of their DRMAA interface implementation. From this documentation they learn how to configure their DRMS in a way that "-l X=true" is used for "X" jobs at site A while "-p 1023" is used at site B for those jobs.

As far as the DRMAA interface specification is concerned only a standardized mechanism for specifying the category is required. The mechanism for associating the policy related portion of the submit command line to the job is to be delivered by each DRMAA implementation. A standardization of this mechanism is

beyond the DRMAA standardization effort, because it is too much related to the administrative interface and it is anticipated that for different DRMS different mechanisms will be appropriate.

2.2.2 Native resource specification

The benefit of the categories concept from the last chapter is that it provides a means for completely hiding site-specific policy details to be considered with a DRMAA job submission for a whole class of jobs. The drawback however of this concept is that it **requires** one job category to be maintained for each policy to be used.

To allow the DRMAA interface to be used also for submission of jobs where job-individual policy specification is required "native resource specification" is supported. Native resource specification can be used without the requirement to maintain job categories. Instead of specifying a category name and having the DRMAA implementation associate the corresponding job submit options, the use of native resource specification will allow directly specifying these submit options.

An example can help to illustrate this idea:

In order to implement the example from section 2.2.1 via native resource specifications, the native option string "-l X=true" had to be passed directly to the DRMAA interface while "-p 1023" had to be used at site B.

As far as the DRMAA interface specification is concerned the native resource specification is an opaque string and interpreted by each DRMAA library. It is possible to use job categories and native resource specification with the same job submission for policy specification. It is assumed that in this case the DRMAA library is capable of joining the outcome of the two policy sources in a reasonable way.

2.3 Interface Routines General Description

The routines are naturally grouped in four categories: init/exit, job submission, job monitoring and control, and auxiliary or system routines like trace file specification and error message routines. All the routines have a prefix "drmaa_".

All of the routines should return an error code upon exit. A possible exception is an auxiliary error message routine that could be modeled after the standard libc strerror routine. In connection to the error routine there should be mechanism for getting DRMAA equivalent of libc errno value. In libc errno is a macro that expands to a modifiable lvalue, such as a dereferenced function pointer to address libc use in reentrant mode.

2.3.1 Init and exit routines

The calling sequence of the init routine should allow all of the considered DRMSs to be properly initialized, either by interfacing to the batch queue commands or to the DRMS API. Likewise, the exit routine should require parameters that will permit proper DRMS disengagement.

2.3.2 Job template and job submission routines

The job submission routines come in two versions. There is one version for submitting individual jobs and one version for submitting bulk jobs. The remote jobs and their attributes are specified with a job template opaque parameter. The job attributes are divided in three groups:

- Core/base or implicit. These have provided setter and getter routines.
- Extended or reserved. These attributes are needed to address desired functionality for particular DRMSs. They are set by using a generic setter routine where the name of the attribute is passed as an extra parameter. All DRMAA libraries need to provide this routine with the default of ignoring the request. Some of these attributes could move to the core/base set.

- Native. These attributes are particular to one or possibly few DRMSs. They are specified via job category mechanism or via the generic setter routine.

The core/base or implicit attributes are:

- Remote command to execute, including the input parameters.
- Job state at submission (suspended/on hold or active).
- Job environment.
- Job working directory.
- Real or wall clock time limit.
- Standard input, output, and error streams.
- E-mail to report the job completion and status.

The extended or reserved attributes are:

- Job execution mode, synchronous or asynchronous.
- Input/output files to be staged and a parameter denoting shared or distributed file system. DRMAA specification 1.0 assumes shared file system. This is to be used with care.
- Job name to be used for the job submission. (Alphanumeric and _ character allowed.)
- Time of execution.

2.3.3 Job monitoring and controlling routines

Job monitoring and controlling API group needs to handle:

- job stopping, resuming, and killing
- waiting for the remote job till the end of its execution
- checking the exit code of the finished remote job
- checking the remote job status
- waiting for all the jobs to finish execution (this is a useful synchronization mechanism)

The Unix and Windows signals are replaced with the job control routines that have counterparts in DRMSs. The only nontraditional feature is the passing of a NULL job_id to indicate operations on all job_ids in the current process.

The remote job could be in following states:

- queued
- system suspended
- user suspended
- running
- finished (un)successfully

To this list we need to add a possibility of DRMAA library not being able to determine the status of the remote job.

2.3.4 Auxiliary routines

The auxiliary routines are needed for execution tracing and error monitoring. The tracing is especially useful for the situations when there is multiple processes spawned few levels deep. The error codes routines and variable drmaa_errno are libc equivalents. drmaa_errno is a macro such as a function reference for reentrant DRMAA library implementation.

The next Chapter contains the an example of an API as specified here.

3 API Specification

The API is preceded by the common constants that are used in the course of the distributed program implementation. For convenience, the API is divided in its four logical sections: init/exit, job submission, job monitoring and control, and auxiliary routines.

To prevent excessive number of job template setter/getter attribute routines, an alternative approach is given in section 3.2.

Disclaimer #1: The code is used here for illustrative purposes. It is not meant to be an example of a full solution at this stage.

Disclaimer #2: The routine names are tentative.

3.1 C/C++ API, explicit job template setter/getter attribute routines

```
/* DRMAA constants */
#define DRMAA_JOB_SIZE 128 /* the size of the job_id buffer */

/* remote job status constants, PS stands for process status */
/* constant that indicates that remote job status cannot be determined */
#define DRMAA_PS_UNDETERMINED -6
/* constant that indicates that remote job is queued */
#define DRMAA_PS_QUEUED -4
/* constant that indicates that remote job is system suspended */
#define DRMAA_PS_SYSTEM_SUSPENDED -3
/* constant that indicates that remote job is user suspended */
#define DRMAA_PS_USER_SUSPENDED -2
/* constant that indicates that remote job is running */
#define DRMAA_PS_RUNNING -1
/* constant that indicates that remote job finished normally */
#define DRMAA_PS_DONE 0
/* constant that indicates that remote job finished, but failed */
#define DRMAA_PS_FAILED 1

/* enumeration for job control specifying */
typedef enum drmaa_control {
    suspend = 0,
    resume,
    terminate,
} drmaa_control_t;

external struct drmaa_job_template;

/* ----- init/exit routines ----- */

/* init DRMAA/DRMS library */
int drmaa_init(char *hostname, int hostport, char *program_name);

/* Disengage from DRMAA/DRMS library and clean up the objects we created */
int drmaa_exit( void );

/* ----- job template and job submission routines ----- */

/* allocate job_template */
drmaa_job_template* drmaa_allocate_job_template( void );
```

```

/* delete job template */
void drmaa_delete_job_template( drmaa_job_template *jt );

/* set a generic attribute */
int drmaa_set_generic_attribute( drmaa_job_template *jt, char *attr_name, char *value );

/* set an attrName attribute, the routines are replicated */
int drmaa_set_attrName( drmaa_job_template *jt, char *value );

/* get attribute value for attribute attrName */
char* drmaa_get_attrName( drmaa_job_template *jt );

/* run a remote job, job_id space allocated by the user */
int drmaa_run_job( char *job_id, [char *job_name,] drmaa_job_template *jt );

/* run parametrized remote jobs */
int drmaa_run_bulk_job( char **job_id, [char *job_name,] drmaa_job_template *jt, int njobs );

/* ----- job control routines ----- */

/* (re)start/stop/kill the job, all if NULL, synchronous used for job (re)starting only */
int drmaa_control( char *job_id, drmaa_control_t mode, int synchronous );

/* synchronize or block till these jobs, all if NULL, of my submitted jobs have finished execution */
int drmaa_synchronize( char **jobs_ids );

/* wait for the remote job to finish and get the remote job exit code, modeled after wait3(2) system call */
int drmaa_waitpid( char *job_id, int *exit_code, int options, char **rusage );

/* get program status of the remote job */
int drmaa_job_ps( char *job_id, int *remote_ps );

/* get the exit code of the finished remote job */
int drmaa_getpid_status( char *job_id, int *exit_code );

/* ----- auxiliary routines ----- */

/* specify a file for tracing */
int drmaa_set_trace_file( char *file_name );

/* pass the text to output in the trace file or stderr by default */
int drmaa_trace_text( char *text );

/* record the error in the trace file, or stderr by default */
int drmaa_perror( char *text );

/* get the error message for the error number */
char *drmaa_strerror( int error );

```

3.2 *Enumerated parameter job template setter/getter attribute routine*

The idea here is to prevent having proliferation of number of setter and getter routines. Instead, `drmaa_basic_attributes_t` enumeration has all the basic/core attributes listed. `drmaa_set_attrName` and `drmaa_get_attrName` are only one affected from the previous section. `drmaa_basic_attributes_t` enumeration and their replacements are defined as:

```
typedef enum drmaa_basic_attributes {  
    command = 0,  
    email,  
    stdinput,  
    stdoutput,  
    stderr,  
    time_limit,  
    initial_job_state,  
} drmaa_basic_attributes_t;  
  
/* set an attrName attribute */  
int drmaa_set_attribute( drmaa_job_template *jt, drmaa_basic_attributes_t dbat, char *value );  
  
/* get attribute value for attribute attrName */  
char* drmaa_get_attribute( drmaa_job_template *jt, drmaa_basic_attributes_t dbat);
```