# UNICORE TARGET SYSTEM INTERFACE

**Sven van den Berghe,** *Fujitsu Laboratories of Europe*

**July 04, GetDirectory replaced by GetFileChunk, FREEZE command**

**27<sup>th</sup> March 03, clarified postcondition on EndProcessing**

**Version 1.0.6, July 04**

.

## 1       General

This document describes the API to the TSI as used by the NJS. This API can be used to develop TSIs for new systems. The parts of the TSI that interact with the target system have been isolated and are documented here with their function calls.

Note that this document is not a complete definition of the API, it is a general overview. The full API specification can be derived by reading the TSI code supplied with a Unicore release.

The functions are implemented in the TSI as calls to Perl methods (with the methods loaded through modules).Input data from the NJS is passed as arguments to the method. Output is returned to the NJS by calling some global methods documented below or by directly accessing the TSI's command and data channels.

TSIs will be shipped with complete implementations of all the functions and can be tailored by changing the supplied code or by implementing new versions of the functions that need to change for the system.

### 1.1     Multithreading

The TSI implementation is single threaded. However, the NJS has been designed to be able to use multiple TSIs. This is done by having more than one TSI worker process running and so any replacement code must be able to handle concurrent calls correctly.

## 2       Services provided by the main TSI

### 2.1     Initialisation

The main TSI will contact the NJS and create the necessary communications. It will receive any initialisation information send by the NJS and then process each command from the NJS and call the appropriate method.

### 2.2     Messages to the NJS

The TSI provides methods to pass messages to the NJS.

In particular the NJS expects every method to call either ok_report or failed_report at the end of its execution.

The messaging methods are:

**ok_report(string)**

> Sends a message to the NJS to say that execution of the command was successful. The *string* is also logged as a debug message.

**failed_report(string)**

> Sends a message to the NJS to say that execution of the command failed. The *string* is sent to the NJS as part of the failure message. It is also logged.

**debug_report(string)**

> Logs *string* as a debug message.

## 2.3 User identity and environment setting

In production mode the TSI will be started as a privileged user capable of changing the process' uid and gid to the user and account[1] requested by the Unicore user. This change is made before the TSI executes any external actions.

The TSI performs three types of work: the execution and monitoring of jobs prepared by the user, transfer and manipulation of files on Storage Servers and the management of Uspaces (including spooled files, outcomes and streamed files). Only the first type of work, execution of jobs, needs a complete user environment. The other two types of TSI work use a restricted set of standard commands (mkdir, cp, rm etc) and should not require access to specific environments set up by users. Furthermore, job execution is not done directly by the TSI but is passed off to the local Batch Subsystem which ensures that a full user environment is set before a job is executed. Therefore, the TSI only needs to set a limited user environment for any child processes that it creates.

The TSI sets the following environment in any child process:

$USER. This is set to the user name supplied by the NJS.

$LOGNAME. This is set to the user name supplied by the NJS.

$HOME This is set to the home directory of the user as given by the target system's password file.

$PATH This is inherited from the parent TSI process (see the "tsi" script file)

Localisations of the TSI can also set any other environment necessary to access the BSS. This is done through the Perl ENV array.

## 2.4 General

If the NJS is set up to do resource or QoS checking, then the incarnations of AbstractTasks that claim an issued Ticket will contain the following line

**#TSI_TICKET <ticket>**

Where <ticket> is a string that can be used to claim some resources allocated to this AbstractTask by the BSS or other parts of the local system.

---

[1] In a test environment the TSI may be started as a non-privileged user and so no changing of uid and gid is possible.

## 3 submit

This function submits the script to the BSS.

### 3.1 Input

The script to be executed.

The string from the NJS is processed to replace all instances of $USER by the user's name and $HOME by the user's home directory.

No further processing needs to be done on the script.

The NJS will embed information in the script that the TSI may need to use. This information will be embedded as comments so no further processing is needed.

Each piece of information will be on a separate line with the format:

```
#TSI_name value
```

If the value is the string "NONE", then the particular information should not be supplied to the BSS during submission.

The information is:

### #TSI_JOBNAME

This is the name that should be given to the job.

If this is "NONE", the TSI will use a default jobname.

### #TSI_OUTCOME_DIR

The NJS expects that the stdout and stderr of the job are written to files named stdout and stderr in the directory named as the value of this field.

The actual requirement is that the stdout and stderr files are in the outcome directory when end_processing() tells the NJS that the job has definitely completed execution. The information is provided here because it is possible to tell BSSs where to put stdout and stderr when a job is submitted.

### #TSI_USPACE_DIR

The initial working directory of the script (i.e. the Uspace directory).

### #TSI_TIME

The run time (wall clock) limit requested by this job in seconds

### #TSI_MEMORY

The memory requirement of the job (in megabytes).

The NJS is told through the IDB if this should be supplied as per processor, per node or per job.

### #TSI_PROCESSORS

The number of processors per node required by the job.

### #TSI_NODES

The number of nodes required by this job.

If the value is 0, then this job is a non-parallel job. If the value is 1, then this job has some parallel characteristic.

**#TSI_FASTFS**

The amount of storage to allocate on a fast temporary file system for this job (in megabytes). On the Uspace StorageServer.

**#TSI_LARGEFS**

The amount of storage to allocate on a large temporary file system for this job (in megabytes). On the Uspace StorageServer.

**#TSI_HOMEFASTFS**

The amount of storage to allocate on a fast temporary file system for this job (in megabytes). On the Home StorageServer.

**#TSI_HOMELARGEFS**

The amount of storage to allocate on a large temporary file system for this job (in megabytes). On the Home StorageServer.

**#TSI_STORAGE_REQUEST <directory> <size>**

A request for **size** megabytes of storage in **directory.** This can appear more than once.

**#TSI_QUEUE**

The BSS queue to which this job should be submitted.

**#TSI_EMAIL**

The email address to which the BSS should send any status change emails.

**#TSI_SWR<name>**

A SoftwareResource requested by the job. These names are site specific. One entry for each Software Resource requested by the job.

**#TSI_INFO <tag> <value>**

An InformationResource contained within the resources requested by the job.

**#TSI_PREFER_INTERACTIVE <junk>**

The presence of this indicates that the task contained a SoftwareResource whose invocation definition in the IDB says that it should be executed "interactively". The stdout and stderr files should be placed in the Outcome directory. However, the TSI can reply with an OK and not the BSS id.

## 3.2    Output

### 3.2.1    *Normal*

**Post condition:** the script is a job on the BSS

The output is the BSS identifier of the job (also used in abort_job, cancel_job, hold_job, resume_job, get_status_listing and end_processing) unless the execution was interactive (in this case the execution is complete when the TSI returns from this call and the output is that from ok_report()).

### 3.2.2    *Error*

failed_report() called with the reason for failure

# 4    get_file_chunk

This function is called by the NJS to fetch the contents of a file.

## 4.1    Input

Introduction:

```
#TSI_FILECHUNK
```

The full path name of the file whose contents need to be sent to the NJS:

```
#TSI_FILE <file name>
```

Where to start reading the file:

```
#TSI_START <start byte – integer>
```

How many bytes to retsurn:

```
#TSI_LENGTH <chunk length – integer>
```

The file name is modified by the TSI to substitute all occurrences of the string "$USER" by the name of the user and all occurrences of the string "$HOME" by the home directory of the user.

## 4.2    Output

### 4.2.1    *Normal*

**Post conditions**

The NJS has a copy of the request part of the file.

### 4.2.2    *Error*

failed_report() called with the reason for failure.


# 5    execute_script

This function executes the script directly from the TSI process, without submitting the script to the batch subsystem. This function is used by the NJS to manipulate the Uspace and Portfolios and to perform some file management functions.

The NJS does not use this to execute any user scripts.

## 5.1    Input

The script to be executed.

The string from the NJS is processed to replace all instances of $USER by the user's name and $HOME by the user's home directory.

No further processing needs to be done on the script.

## 5.2    Output

### 5.2.1    *Normal*

**Post condition:** The script has been executed

Concatenated stderr and stdout from the execution of the script is sent to the NJS following the  ok_report() call.

### 5.2.2    *Error*

failed_report() called with the reason for failure.

# 6    abort_job

This function sends a command to the BSS to abort the named BSS job. Any stdout and stderr produced by the job before the abort takes effect must be saved.

The NJS will follow this call with a call to end_processing.

## 6.1    Input

The BSS identifier of the job to abort as the string *identifier* in:

```
#TSI_BSSID identifier
```

## 6.2    Output

### 6.2.1    *Normal*

**Post condition:** The job is no longer executing on the BSS.

No extra output.

### 6.2.2    *Error*

failed_report() called with the reason for failure.

# 7 cancel_job

This function sends a command to the BSS to cancel the named BSS job. Cancelling means both finishing execution on the BSS (as for abort) and removing any stdout and stderr.

The NJS will **not** follow this call with a call to end_processing.

## 7.1 Input

The BSS identifier of the job to cancel as the string *identifier* in:

```
#TSI_BSSID identifier
```

## 7.2 Output

### 7.2.1 *Normal*

**Post condition:** The job is no longer execution and stdout and stderr have been deleted.

No extra output.

### 7.2.2 *Error*

failed_report() called with the reason for failure.

## 8    hold_job

This function sends a command to the BSS to hold execution of the named BSS job. Holding means suspending execution of a job that has started or not starting execution of a queued job.

Note that suspending execution can result in the resources allocated to the job being held by the job even though it is not executing and so some sites may not allow this. This is dealt with by the relaxed post condition below[2].

Some sites can hold a job's execution and release the resources held by the job (leaving the job on the BSS so that it can resume execution). This is called freezing. The NJS can send a request for a freeze which the TSI may execute, if there is no freeze command initialised the TSI may execute a hold in its place

### 8.1    Input

Request a freeze if possible:

```
#TSI_FREEZE
```

The BSS identifier of the job to hold as the string *identifier* in:

```
#TSI_BSSID identifier
```

### 8.2    Output

#### 8.2.1  *Normal*

**Post condition:** true.

No extra output.

#### 8.2.2  *Error*

failed_report() called with the reason for failure.

---

[2] So an acceptable implementation is for hold_job to return without executing a command.

# 9 resume_job

This function sends a command to the BSS to resume execution of the named BSS job.

Not that suspending execution can result in the resources allocated to the job being held by the job even though it is not executing and so some sites may not allow this. This is dealt with by the relaxed post condition below[3].

## 9.1 Input

The BSS identifier of the job to resume as the string *identifier* in:

```
#TSI_BSSID identifier
```

## 9.2 Output

### 9.2.1 Normal

**Post condition:** the job is executing on the BSS or is queued and will start executing when its turn comes.

No extra output.

### 9.2.2 Error

failed_report() called with the reason for failure.

---

[3] An acceptable implementation is for resume_job to return without executing a command (if hold_job did the same).

# 10      get_status_listing

This function returns the status of all the jobs on the BSS that have been submitted through any TSI providing access to the BSS.

## 10.1    Input

None.

This method is called with the TSI's identity set to the special QSTAT_XLOGIN user from the NJS. This is because the NJS expects the returned listing to contain every Unicore job from every Unicore user but some BSS only allow a view of the status of all jobs to privileged users.

## 10.2    Output

### 10.2.1  *Normal*

**Post condition:** the NJS has an up to date list of the state of the jobs on the BSS.

The first line is "QSTAT\n".

There follows an arbitrary number of lines, each line containing the status of a job on the BSS with the following format:

```
identifier status
```

Where *identifier* is the BSS identifier of the job and *status* is one of: QUEUED, RUNNING or SUSPENDED.

The output **must** include all jobs still on the BSS that were submitted by a TSI executing on the target system (including all those submitted by TSIs other than the one executing this command). The output **may** include lines for jobs on the BSS submitted by other means (the NJS will ignore these lines).

### 10.2.2  *Error*

failed_report() called with the reason for failure.

# 11      end_processing

This function is called by the NJS when it suspects that a job executing on the BSS has completed. This function confirms that the job has finished executing and returns part of the stderr produced by the job.

After successfully executing this function (on a completed job) stdout and stderr must be in the outcome directory (unless the disposition is to delete the files).

The files can be moved into the Outcome directory by this function or by the way that the BSS submit command was used.

## 11.1      Input

The BSS identifier of the job to check as the string *identifier* in:

        TSI_BSSID *identifier.*

The directory for the result files in:

        TSI_OUTCOME_DIR *directory_name*

The *directory_name* is modified by the TSI to substitute all occurrences of the string "$USER" by the name of the user and all occurrences of the string "$HOME" by the home directory of the user.

The final disposition of the files in:

        TSI_DISPOSITION *option*

Where *option* is either KEEP or DELETE

## 11.2      Output

### 11.2.1      *Normal*

**Post conditions:**

*Either*:

There exist files named **stdout** and **stderr**  in the AbstractAction's Outcome directory and the stderr file has a line containing the string "UNICORE EXIT STATUS". The NJS is sent the whole of this line. The NJS is also sent the whole line (if any) containing the string "UNICORE DECISION".

The NJS interprets this as a completed batch job.

*Or:*

Neither **stdout** or **stderr** exist. The NJS is sent "TSI_STILLEXECUTING" with a call of ok_report().

The NJS interprets this as an incomplete batch job (transiently disappeared from queue).

*Or:*

The files **stdout** and **stderr** exist but stderr does not contain the string "UNICORE EXIT STATUS". No special message is sent to the NJS (backwards compatibility).

NJS interprets this as an incomplete batch job (stdout/stderr still being written)

Note that if for some reason the exit status will never be written the NJS will never mark the batch job as complete. The user will have to detect this manually and abort or cancel the AbstractAction (or reply on a cleanup from TerminationTime).

### 11.2.2 *Error*

failed_report() called with the reason for failure.

## 12    put_files

This function is called by the NJS to write the contents of one or more files to a directory accessible by the TSI.

### 12.1    Input

TSI_FILESACTION contains the action to take if the file exists (or does not):

- 0 = don't care
- 1 = only write if the file does **not** exist
- 2 = only write if the file exists

This action applies to all the files is a call of put_files.

The files are read from the data channel following this pseudo code:

```
while(files to transfer) {

   read filename and permissions from command channel

   substitute all occurrences of the string "$USER" by
the name of the user and all occurrences of the string
"$HOME" by the home directory of the user.

   while(file has more bytes) {

      read packet_size from command channel

      read packet_size bytes from data channel

      write bytes to file

   }

}
```

Where "permissions" are the owner's permissions to set on the file (following the convention of the previous section).

### 12.2    Output

*12.2.1    Normal*

**Post conditions:** The TSI has written the files to the directory.

None

*12.2.2    Error*

failed_report() called with the reason for failure.

## 13    Version History

**8<sup>th</sup> January 03** Changes for Unicore 4

**31 October 01** Extra information in a PutFiles call

**28 September 01** Uspace passed to ExecuteScript, interactive flag to submit

**23 July 01** Added the file owner's permissions for GetDirectory and PutFiles

## 14    PUBLIC LICENSE    (VERSION 1.0)

The Software (as defined below) is made available under the terms of this Public License Agreement   Any use, installation, reproduction, modification or distribution of the Software by You constitutes Your acceptance of the terms contained herein.

### 14.1   Definitions

 "Contributor" means a person or an organization which has made changes or additions to the Software and has authorized the use, installation, reproduction, modification and distribution of its changes or additions with the Software in accordance with the terms of this License.  Contributor includes the original copyright holders of the Original Code.

 "Executable Code" means the Software in any form other than Source Code.

 "Larger Work" means a work which combines Original Code or portions thereof with code not governed by the terms of this License.

"License" means this Public License Agreement.

"Original Code" means the Source Code of computer software as originally made available under this License and  which is described in header file of such software or in the Source Code notice in Exhibit A as Original Code.

"Software" means the Executable Code, the Original Code or Modifications or the combination of the Original Code and Modifications, in each case including portions thereof.  Software may include voluntary contributions made by Contributors.

"Source Code" means the preferred form of the Software for making modifications to it, including all modules it contains, plus any associated documentation, interface definition files, scripts used to control compilation and installation of an Executable or source code differential comparisons against either the Original Software or another well known, available Covered Software of the Contributor's choice.  The Source Code can be in a compressed or archival form, provided the appropriate decompression or de-archiving software is widely available for no charge.

"Modifications" means any addition to or deletion from the substance or structure of either the Original Code or any previous Modifications.  When Software is released as a series of files a Modification is any addition to or deletion from the contents of a file containing any part of the Original Code or previous Modifications or any new file that contains any part of the Original Code or previous Modifications.

 "You" or "Your" means an individual or a legal entity exercising rights under and complying with all of the terms of this License or a future version of this License.

## 14.2   License Terms

Installation, use, reproduction, display, Modification and redistribution of this Software, with  or without Modification, in Source and Executable forms are permitted by You and Your sub-licensees under this License subject to the following conditions:

1. All copies and redistributions of this Software, with or without Modification, must reproduce the terms of this License in: (1) the Software, and (2) the user documentation or some other similar material which is provided with the Software (if any).

2. The user documentation, if any, included with redistribution, must include the following notice:

   "This product includes software developed by Fujitsu Limited (http://www.fujitsu.com)."

   This acknowledgment notice must also be reproduced in the Software itself if that is where third-party acknowledgments normally appear.

3. In the event that any trademark(s) belonging to Fujitsu Limited appear on or are associated with the Software such trademark(s) may not be used to endorse or promote software, or products derived therefrom, and may not be affixed to modified redistributions of this Software except with prior written approval, obtainable from Fujitsu Limited.  Except as expressly stated in this License You receive no rights or licenses to the intellectual property of any Contributor under this License whether expressly, by implication, estoppel or otherwise.  All rights in the Software not expressly granted under this Agreement are reserved.

4. To the extent that patent claims licensable by the Fujitsu Limited are necessarily infringed by the use or redistribution of the Software, You are granted a non-exclusive, worldwide, royalty-free perpetual license under such patent claims, with the rights for You to make, use, license, offer to license, import and otherwise transfer the Software in Source Code and Executable Code form and Larger Works.   This patent license shall apply to the combination of the Software with other software if, at the time the Software is added by You,  such addition of the Software causes such combination to be covered by such patent claims. This patent license shall not apply to any other combinations which include the Software.  No hardware per se is licensed hereunder.

5. If You or any subsequent sub-licensee (a "Recipient") institutes patent litigation against any Contributor (including a cross-claim or counterclaim in a lawsuit) alleging that the Software infringes such Recipient's patent(s), then Your rights or such Recipient's rights granted (directly or indirectly) under the patent license above shall terminate as of the date such litigation is filed.  All sublicenses to the Software which have been properly granted prior to termination shall survive any termination of said patent license, if not otherwise terminated pursuant to this section.

6. THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY, OF SATISFACTORY QUALITY AND FITNESS FOR A PARTICULAR PURPOSE OR USE ARE DISCLAIMED.  By Your acceptance of this License, You understand that although the copyright holders and the Contributors each grant the licenses to its contributions, no assurances are provided to You by any  copyright or intellectual

property right holder that the Software and/or any contribution(s) thereto does not infringe the patent or other intellectual property rights of any other entity. The original copyright holders and Contributors disclaim any liability to You and to any Recipient for claims brought by any other entity based on infringement of intellectual property  rights or otherwise.  As a condition to  exercising the rights and licenses granted hereunder You hereby assume sole responsibility to secure any other intellectual property rights needed, if any.

7.  You are solely responsible for determining the appropriateness of using and distributing the Software and You assume all risks associated with the exercise of the rights granted to you under this License,  including but not limited to the risks and costs of program errors, compliance with applicable laws, damage to or loss or data, programs or equipment and unavailability or interruption of operations.

8.  THE COPYRIGHT HOLDERS AND CONTRIBUTORS SHALL HAVE NO LIABILITY TO YOU, TO ANY RECIPIENT OR TO ANY OTHER PERSONS FOR DIRECT, INDIRECT, SPECIAL, INCIDENTAL, CONSEQUENTIAL, EXEMPLARY, OR PUNITIVE DAMAGES OF ANY CHARACTER INCLUDING, WITHOUT LIMITATION, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES, LOSS OF USE, DATA OR PROFITS, OR BUSINESS INTERRUPTION, HOWSOEVER CAUSED AND ON ANY THEORY OF CONTRACT, WARRANTY, TORT (INCLUDING NEGLIGENCE), PRODUCT LIABILITY OR OTHERWISE, ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

9.  Commercial distributors of software may accept certain responsibilities with respect to end users, business partners and the like.  While this License is intended to facilitate the commercial use of the Software, if You include the Software in a commercial product offering You should do so in a manner which does not create potential liability for the copyright holders or any other Contributors.  If You include the Software in a commercial product offering then You hereby agree to defend and indemnify the copyright holders and the Contributors ("Indemnified Contributor") against any losses, damages and costs (collectively "Losses") arising from claims, lawsuits and other legal actions brought by a third party against the Indemnified Contributor to the extent caused by the acts or omissions of You in connection with Your distribution of the Software in a commercial product offering.  The obligations in this section do not apply to any claims or Losses relating to any actual or alleged intellectual property infringement.  In order to qualify, an Indemnified Contributor must: a) promptly notify You in writing of such claim, and b) allow You to control, and cooperate with You in, the defense and any related settlement negotiations.  The Indemnified Contributor may participate in any such claim at its own expense.

10. All of Your rights under this License shall terminate automatically if You fail to comply with any of the terms or conditions of this License and do not cure such failure within a reasonable period of time after becoming aware of such non-compliance.  If Your rights under this License terminate, You agree to cease use and distribution of the Software as soon as reasonably practicable.  Your obligations under this License and any sub-licenses granted by You relating to the Software hereunder prior to termination shall continue and survive.

11. Fujitsu Limited may publish revised and/or new versions of this License from time to time.  Each version will be given a distinguishing version number.  Once Software has been published under a particular version of this License, You may

continue to use it under the terms of that version. You may also choose to use such Software under the terms of any subsequent version of this License published by Fujitsu Limited.  No one other than Fujitsu Limited has the right to modify the terms of this License.

12. This Public License Agreement is governed by the laws of Japan.   No party to this License will bring a legal action under this Agreement more than one year after the cause of action arose.