

Getting Started with the Serpens extension for Kepler-2.3

Working with vine:toolkit module

{marcinp, michalo, tzok}@man.poznan.pl

vine:toolkit module description covers:

Data management using RAS/vine:toolkit

Job management using RAS/vine:toolkit

Example workflows

Template workflow for job submission

This paper covers Serpens 2.4 release

January 30, 2012

Getting Started with the Serpens - vine:toolkit module

The *Getting Started with the Serpens - vine:toolkit module* guide is for scientists who would like to use the Serpens suite in Kepler.

The Serpens was created to assist scientists who would like to use Kepler to manage GRID based job execution from the level of Kepler-2.3 workflow.

Table of Contents

1. Introduction.....	2
1.1 What is vine:toolkit?.....	3
1.2 What is the Serpens - vine:toolkit module?.....	3
1.3 Accessing HPC resources via RAS/Vine.....	4
2. Using the Serpens - vine:toolkit module.....	5
2.1 Authentication/authorization within Serpens vine:toolkit. .	5
2.2 Uploading files into UNICORE storage via vine:toolkit.....	6
2.2.1 RASUploadFile - example workflow.....	7
2.3 Downloading files from UNICORE storage.....	7
2.3.1 RASDownloadFile - example workflow.....	8
2.4 Registering output files at UNICORE storage.....	8
2.5 Creating JSDL based description of the job - RASJSDLGenerator.....	10
2.5.1 RASJSDLGenerator - example workflow.....	11
2.6 Submitting jobs using RASSubmit actor.....	11
2.6.1 RASSubmit - example workflow.....	12
2.7 Checking status of the job using RASStatusLoop actor....	12
2.7.1 RASStatusLoop - example workflow.....	13
2.7 Retrieving job's output/error streams.....	13
2.8 Putting things together.....	14
References.....	15
Acknowledgments.....	15

1. Introduction

The *Getting Started with the Serpens - vine:toolkit module* guide introduces the main components and functionality of the *Serpens - vine:toolkit* module. The *Serpens - vine:toolkit module* is an add-on module suite for Kepler, a software application for the analysis and modeling of scientific data, and it provides functionality for managing data management and job management within GRID environment from the level of Kepler workflow.

1.1 What is vine:toolkit?

From official website <http://vinetoolkit.org>:

Vine is a modular, extensible Java library that offers developers an easy-to-use, high-level Application Programmer Interface (API) for Grid-enabling applications. Vine can be deployed for use in desktop, Java Web Start, Java Servlet 2.3 and Java Portlet 1.0 environments with ease. Plus Vine supports a wide array of middleware and third-party services, so you can focus on your applications and not lose focus on the Grid!

Supported Middleware and Standards:

- QosCosGrid
- BES
- gLite3
- Globus Toolkit 4.0.x, 4.2.1
- GRIA 5.3
- Unicore 6
- JSDL
- Storage Resource Broker
- Storage Resource Manager
- RUS
- OGSA-DAI 2.2
- Active Directory

1.2 What is the Serpens - vine:toolkit module?

The *Serpens - vine:toolkit module* provides a way for accessing GRID resources through the Kepler workflow.

The *Serpens - vine:toolkit module* allow users to:

- upload/download data to/from UNICORE storage
- submit jobs into UNICORE
- check statuses of the running jobs
- retrieve *Standard Output* and *Standard Error* of HPC job
- retrieve HPC job outcome

Serpens - vine:toolkit module is based on vine:toolkit, thus it is heavily based on so called Domain concept. You can read more about Domains at vine:toolkit home page¹. Simply saying - domain name describes resources

1 <http://vinetoolkit.org/content/domainxml-multi-domain-configuration>

available for user. Within this document DN or simply domain will refer to Domain Name. It is required by user to know names of domains available at given Vine:toolkit installation.

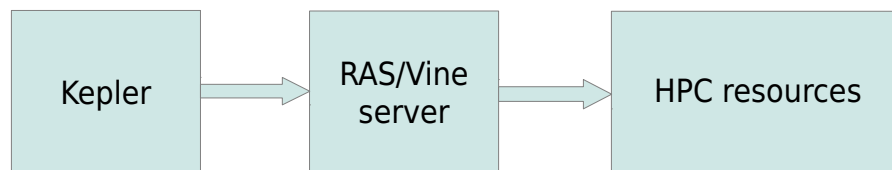
Important, Serpens vine:toolkit module was tested with UNICORE versions up to 6.3.1.

1.3 Accessing HPC resources via RAS/Vine

Serpens - vine:toolkit module is based on RAS - Roaming Access Server. RAS is a middleware that allows external clients to access HPC/GRID resources. Access to resources is granted basing on user's GRID certificates (X.509).

Each actor performing tasks related to HPC access uses RAS to achieve the task.

Architecture of the solution follows:



This architecture allows users to benefit following ways:

1. HPC resources are exposed via RAS/Vine - this way configuration takes place only at server side
2. RAS/Vine provides access to other infrastructures - this way client can benefit from universal GRID/HPC access layer
3. RAS/Vine is able to expose resources that are not available due to security related concerns (e.g. are hidden behind firewall)
4. RAS/Vine allows clients to access resources from behind the firewall (e.g. client is using network from highly secured environment - in that case there is a need for only one, http, port to be opened).

Most of the workflows presented in this documentation use so called **hostAddressRAS**. This is the address of RAS/Vine server used for communication between client (Kepler) and resources (HPC).

2. Using the Serpens - vine:toolkit module

Once you have installed the Serpens - vine:toolkit module, you will notice new actors available at the Component list.

Actor name	Purpose
RASUploadFile	Uploads file into UNICORE storage via vine:toolkit
RASDownloadFile	Downloads file from the UNICORE storage via vine:toolkit
RASRegisterFile	Registers file at UNICORE storage via vine:toolkit (required for output files)
RASGetError	Retrieves job's error stream
RASGetOutput	Retrieves job's output stream
RASJSDLGenerator	Generates JSDL description of the job basing on input parameters
RASStatusLoop	Checks job's status
RASSubmit	Submits job into UNICORE via vine:toolkit
RASUName	Converts file name relative to job into UNICORE description of file
RASGenerateProxy	Creates user's proxy basing on X.509 certificate
RASVomsProxy	Creates user's VOMS proxy basing on user's proxy.

2.1 Authentication/authorization within Serpens vine:toolkit

Access to HPC resources is restricted and requires authentication/authorization to reach certain resources. Some actors require user to be properly authenticated before an action can take place. These actors are: *RASUploadFile*, *RASDownloadFile*, *RASRegisterFile*, *RASSubmit*, *RASStatusLoop*, *RASGetError*, *RASGetOutput*.

Serpens vine:toolkit authentication is based on VOMS² proxy concept. VOMS proxy can be generated using two actors:

- *RASGenerateProxy* - generates proxy basing on user's X.509 certificate
- *RASVomsProxy* - generates VOMS proxy basing on the outcome of *RASGenerateProxy* actor

Example usage of these actors is shown on following picture:



workflow type: GRID/HPC

In order to get applications running using Serpens you might need to generate GRID proxy or VOMS proxy. In case you need GRID proxy it can be generated through using following workflow.

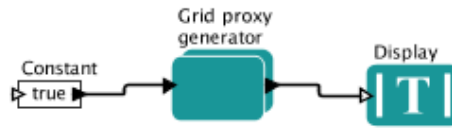


Illustration 1: Generating grid proxy basing on X.509 certificate

This workflow can be found at following location after Serpens demo workflows are installed: `$SERPENS_MODULE/vine-toolkit/demos/gridproxygenerator.xml`³. Please note that you will need your X.509 certificate (both certificate and private key) to be located at your \$HOME directory. Typically certificate and key files will have following names `usercert.pem` and `userkey.pem` respectively.

2.2 Uploading files into UNICORE storage via vine:toolkit

To use external data for processing one must upload data into UNICORE storage to use it during HPC job execution. This can be done using RASUploadFile actor. RASUploadFile actor allows uploading local input files into UNICORE storage specified within vine:toolkit's domain.

NOTE: RASUploadFile/RASDownloadFile actors are much less efficient than UnicoreUploadFile/UnicoreDownloadFile actors. However, they can operate within restricted environment where Internet traffic goes through firewall.

To upload a file into UNICORE storage (using RASUploadFile actor) it is required to provide following information:

Input port name	Description
proxy	String representing user's proxy
domainName	Name of the Domain within vine:toolkit. This parameter allows to change target machine transparently by redirecting vine:toolkit calls to different target registry.
fileName	Local file name. This file will be uploaded into UNICORE storage. Base name will be used as file name at storage.

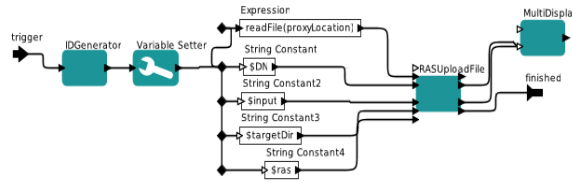
³ Throughout the document it is assumed that \$SERPENS_MODULE points to directory where Serpens module was installed. This is the location where you have checked out Kepler 2.0 and Serpens related modules.

targetDirectory	Target directory at UNICORE storage. If not specified, root directory will be used.
hostAddressRAS	URL of RAS.

2.2.1 RASUploadFile - example workflow

Following example shows example usage of RASUploadFile actor. This workflow will perform following action:

- it will read file from \$SERPENS_MODULE/vine-toolkit/demos/input.txt
- it will upload a file into UNICORE storage
- it will display file location at UNICORE storage



PARAMETERS

<p style="color: red; text-align: center;">Input specification</p> <div style="background-color: #f08080; border-radius: 15px; padding: 5px; text-align: center; color: white; font-weight: bold; margin-bottom: 5px;">INPUT FILES</div> <p style="color: blue; font-size: small;">LOCAL INPUT FILES</p> <ul style="list-style-type: none"> ● input: \$demoLocation/input.txt 	<p style="color: yellow; text-align: center;">You MUST check if these are correct</p> <div style="background-color: #f0f080; border-radius: 15px; padding: 5px; margin-bottom: 5px;"> <ul style="list-style-type: none"> - 'demos/' directory location - proxy file location - RAS address </div> <div style="background-color: #f0f080; border-radius: 15px; padding: 5px; margin-bottom: 5px;"> <ul style="list-style-type: none"> ● proxyLocation: \$demoLocation/proxy ● outputDir: \$demoLocation/output/ </div> <div style="background-color: #f0f080; border-radius: 15px; padding: 5px; margin-bottom: 5px;"> <p style="font-size: x-small; margin: 0;">Target domain defined within RAS</p> <ul style="list-style-type: none"> ● DN_S: seagrass.man.poznan.pl ● DN: \$DN_S </div>	<p style="color: green; text-align: center;">Leave these unchanged</p> <div style="background-color: #90ee90; border-radius: 15px; padding: 5px; margin-bottom: 5px;"> <ul style="list-style-type: none"> - predefined RAS addresses </div> <div style="background-color: #90ee90; border-radius: 15px; padding: 5px; margin-bottom: 5px;"> <ul style="list-style-type: none"> ● ras: http://senecio.man.poznan.pl:8090 </div> <div style="background-color: #90ee90; border-radius: 15px; padding: 5px; margin-bottom: 5px;"> <p style="font-size: x-small; margin: 0;">Location of the files at the storage element. This variables are used internally by workflow.</p> <ul style="list-style-type: none"> ● targetDir: \$uid ● uuid: 604e5a7-e744-4740-a1bf-9ce2b559dc61 ● moduleVersion: 2.2.0 ● demoLocation: \$HOME/KeplerData/kepler.modules/vine-toolkit-smoduleVe </div>
---	--	---

Illustration 2: Uploading local file into UNICORE storage - RASUploadFile

You can find this workflow at following location:

`$SERPENS_MODULE/vine-toolkit/demos/vine_upload.xml`

2.3 Downloading files from UNICORE storage

Vine:toolkit module allows downloading files from UNICORE storage. RASDownloadFile actor allows user to specify file location at UNICORE storage and download file into local directory. Note, that file location at UNICORE storage is internally stored by vine:toolkit for each domain and each user. Thus, users can operate on files with virtual names. Instead of specifying file's URI one can provide only directory name and file name at UNICORE storage (e.g. /my_data/input_file.txt).

To retrieve file from UNICORE storage, following parameters must be specified:

Input port name	Description
proxy	String representing user's proxy

domainName	Name of the Domain within vine:toolkit. This parameter allows to change target machine transparently by redirecting vine:toolkit calls to different target registry.
fileName	Location of file at UNICORE storage.
targetDirectory	Target directory at local file system (download area)
hostAddressRAS	URL of RAS.

2.3.1 RASDownloadFile - example workflow

Following example shows example usage of RASDownloadFile actor. This workflow will perform following action:

- upload a local file into UNICORE storage
- display file address at UNICORE storage
- download a file from UNICORE storage
- store it within local directory
- display local, downloaded file name

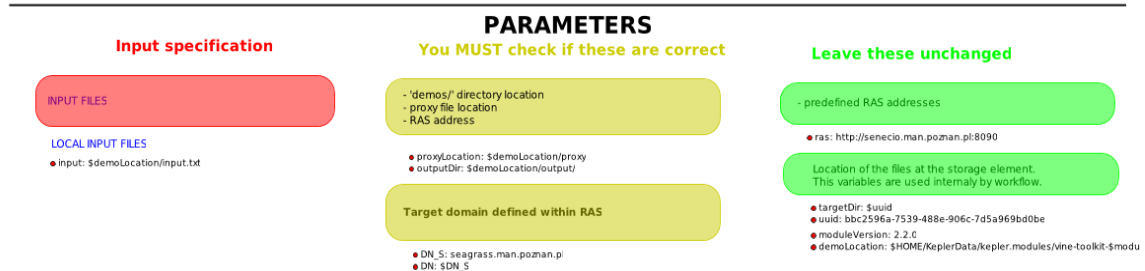
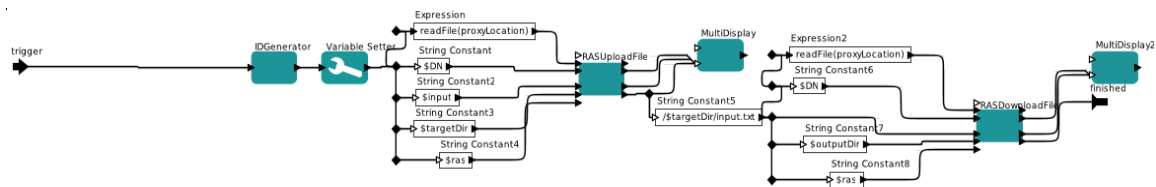


Illustration 3: Downloading file using RASDownloadFile actor

Workflow can be found at:

`$SERPENS_MODULE/vine-toolkit/demo/vine_upload_download.xml`

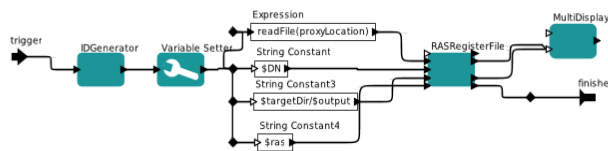
2.4 Registering output files at UNICORE storage

To retrieve files that are generated by HPC job user must specify output location for the files so they can be staged out from job's directory after it is finished. In order to create target location user can register file using RASRegisterFile actor. This actor creates file's URI that describes file location at UNICORE storage. Once created it can be passed to RASJSDLGenerator actor as output file location.

RASRegisterFile actor requires following input parameters:

Input port name	Description
proxy	String representing user's proxy
domainName	Name of the Domain within vine:toolkit. This parameter allows to change target machine transparently by redirecting vine:toolkit calls to different target registry.
fileName	Location of file at UNICORE storage (it can contain parent's names as well, e.g. /my_data/my_job/output_file.txt).
hostAddressRAS	URL of RAS.

As a result it output single String token with URI of the file.



Input specification

INPUT FILES

LOCAL INPUT FILES

- output: output_file.txt

PARAMETERS

You MUST check if these are correct

- 'demos/' directory location
- proxy file location
- RAS address

- proxyLocation: \$demoLocation/proxy
- outputDir: \$demoLocation/output

Target domain defined within RAS

- DN_S: seagrass.man.poznan.pl
- DN: \$DN_S

Leave these unchanged

- predefined RAS addresses

- ras: http://senecio.man.poznan.pl:8090

Location of the files at the storage element.
This variables are used internally by workflow.

- targetDir: \$uuid
- uuid: 5157fa76-abaf-4979-9796-28a5861ef00d
- moduleVersion: 2.2.0
- demoLocation: \$HOME/KeplerData/Kepler.modules/vine-toolkit-\$mod

Illustration 4: Registering output file using RAS/Vine

Workflow can be found at:

`$SERPENS_MODULE/vine-toolkit/demos/vine_register_file.xml`

2.5 Creating JSDL based description of the job - RASJSDLGenerator

To submit a job into UNICORE one must define it via JSDL format. This format allows user to define all the predefined parameters of the job. However, defining JSDL by your own is not convenient - JSDL is an XML subset. In order to provide users with easier way of defining job description Serpens module provides RASJSDLGenerator actor. This actor creates JSDL description of the job basing on input parameters.

Input port name	Description
inputStreamFileName	specifies that given file should be treated as input stream of the application. File must be present in the list of input files
outputStreamFileName	specifies that given file should be treated as output stream of the application. File must be present in the list of output files
inputFileNames	list of file names that will be staged into job's directory
inputFileNamesIURLs	list of input files location at UNICORE storage. NOTE! inputFileNames.length must be equal to inputFileNamesIURLs.length
outputFileNames	list of file names that will be staged out from the job's directory
outputFileNamesURLs	list of output file locations at UNICORE storage. NOTE! outputFileNames.length must be equal to outputFileNamesIURLs.length
executable	executable name at target HPC machine (NOTE. you can either use full path "/bin/ls" or application name when it is system wide available)
arguments	executable arguments
nodes	number of nodes that will be used (NOTE! meaning of this parameter may be HPC machine specific - it might reserve either cores or whole CPU nodes)
individualCPU	number of CPUs reserved for the application (NOTE! meaning of this parameter may be HPC machine specific - it might reserve either cores or whole CPU nodes)
cpuTime	application's wall time

As a result it outputs JSDL description of the job.

2.5.1 RASJSDLGenerator - example workflow

Below there is presented simple RASJSDLGenerator workflow which generates JSDL description of the job. This job is designed to start /bin/lis application with "-la" arguments (vine_jsdl.xml).

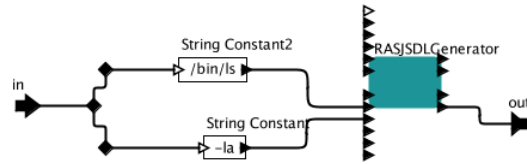


Illustration 5: RASJSDLGenerator example

Workflow can be found at:

`$SERPENS_MODULE/vine-toolkit/demos/vine_jsdl.xml`

2.6 Submitting jobs using RASSubmit actor

To submit a job into UNICORE one must use RASSubmit actor. By providing this actor with input parameters we can define what RAS/Vine host will be used for job submission, which Domain will be used and how job description will look like.

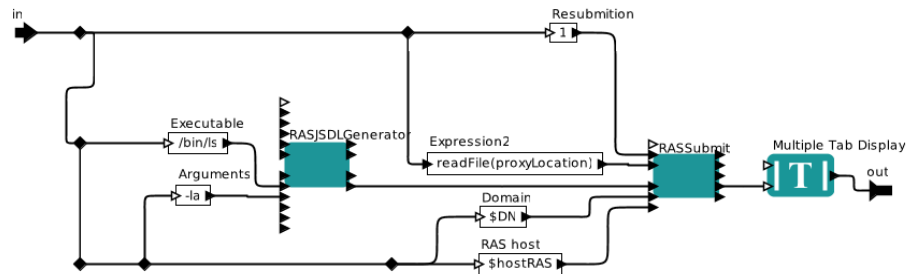
RASSubmit actor accepts following input data:

Input port name	Description
proxy	String representing user's proxy
resubmit	Number of times actor will try to submit the job in case of failure
jsdl	JSDL based description of the job
hostAddressRAS	URL of RAS.
DomainName	Name of the Domain defined at RAS/Vine host

As an output RASSubmit actor generates job id that was submitted via Vine:toolkit library.

2.6.1 RASSubmit - example workflow

Below you can see an example of workflow that submits job after JSDL description is created using RASJSDLGenerator actor (vine_submit.xml).



Target domain defined within RAS and RAS host name

Proxy location

- DN_S: seagrass.man.poznan.pl
- DN: \$DN_S
- hostRAS: http://senecio.man.poznan.pl:8090

- proxyLocation: \$demoLocation/proxy
- moduleVersion: 2.2.0
- demoLocation: \$HOME/KeplerData/kepler.modules/vine-toolkit-\$moduleVersion/demos/

Illustration 6: Submitting jobs using RASSubmit actor

In this particular case JSDL is passed directly from RASJSDLGenerator actor, however it can be read from file as well. It can be also generated by other means.

Workflow can be found at:

`$SERPENS_MODULE/vine-toolkit/demos/vine_submit.xml`

2.7 Checking status of the job using RASStatusLoop actor

RASStatusLoop actor allows user to check the status of previously submitted job. Actor accepts following inputs:

Input port name	Description
proxy	String representing user's proxy
JobId	Id of the job
domainName	Name of the domain that will be used during status check.
hostAddressRAS	URL of RAS.

As a result, actor outputs current status name.

2.7.1 RASStatusLoop - example workflow

Usually, checking the status of the job is executed within the loop. This way, workflow can “wait” till particular job is finished and proceed with the results. You can see example workflow at following image (vine_status.xml):

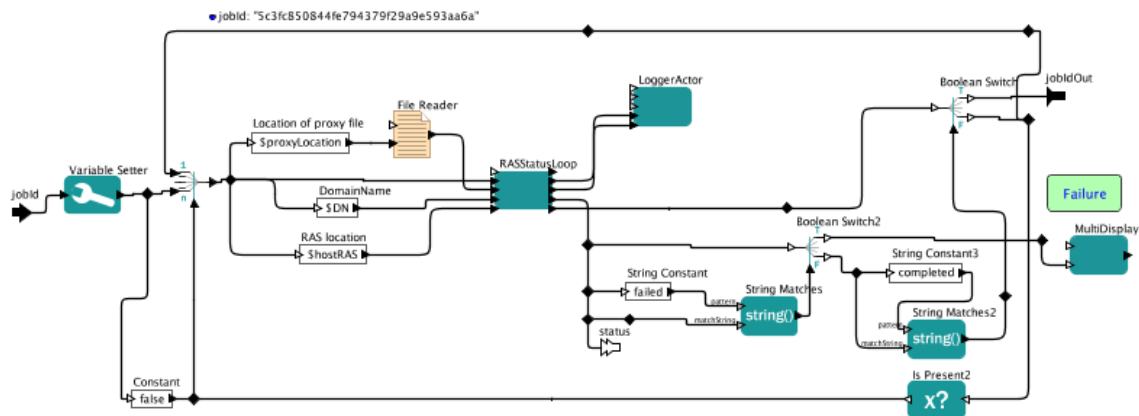


Illustration 7: Checking job's status within loop

Workflow can be found at:

`$SERPENS_MODULE/vine-toolkit/demos/vine_status.xml`

2.7 Retrieving job's output/error streams

To retrieve job's output/error stream user should use RASGetOutput/RASGetError actor. These actors accept following input data:

Input port name	Description
proxy	String representing user's proxy
JobId	Id of the job
domainName	Name of the domain that will be used during status check.
hostAddressRAS	URL of RAS.

They return standard output and standard error respectively.

Following picture presents example usage of these actors:

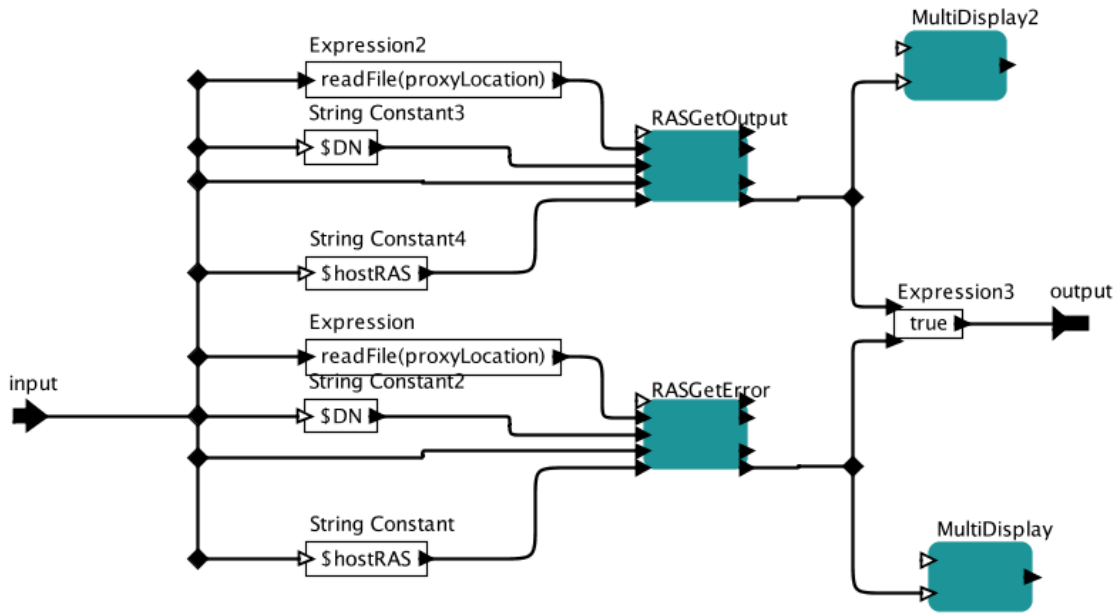


Illustration 8: Retrieving std out/std error of job

This workflow is a part of template workflow described in next section.

2.8 Putting things together

In general, users are not required to use all RAS/Vine related actors in order to build workflow by themselves. Serpens module, apart from actors, provides composite workflows that perform more complicated tasks. An example of such workflow is `vine_template.xml` workflow that allows user to submit job into UNICORE by specifying parameters defined within composite actor. This way, user can simply copy paste one actor in order to get the whole HPC related scenario running - that is: upload files, submit job, check it's status, retrieve outputs. You can see an example of such a workflow at following picture

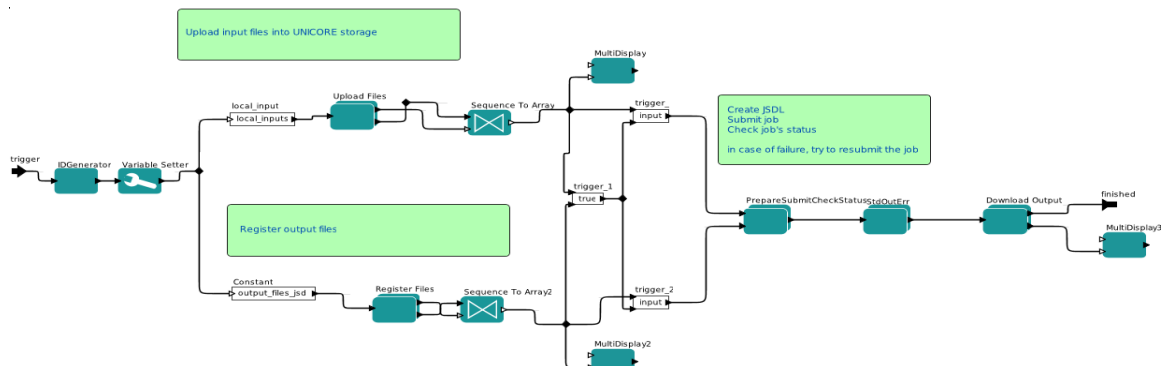


Illustration 9: HPC job submission template

Workflow can be found at:

`$SERPENS_MODULE/vine-toolkit/demos/vine_template.xml`

What you can see on the picture is the workflow that is divided into sections – each section is responsible for an action related to HPC job management. After workflow is finished, output of the job will be stored in specified location.

References

1. Serpens - <http://serpens.psnc.pl/>
2. UNICORE - <http://www.unicore.eu/>
3. Kepler-2.3 - <https://kepler-project.org/>
4. vine:toolkit - <http://vinetoolkit.org/>

Acknowledgments

This research has also received funding from the European Community's Seventh Framework Program (FP7/2007-2013) under grant agreement n°211804 (EUFORIA).