

The Catania Science Gateway for the EGI Long Tail of Science – A How-to

Last modified: 25/11/2016
Version: 1.0

1.	INTRODUCTION.....	1
2.	APPLICATIONS.....	2
2.1	<i>R</i>	2
2.1	<i>Chipster</i>	4
2.3	<i>ClustalW2</i>	6
2.4	<i>Semantic Search Engine</i>	8
3.	MYJOBS	11
4.	SUPPORT	12

1. Introduction

The Catania Science Gateway Framework¹ (CSGF) is an open source toolkit jointly developed by the Division of Catania of the Italian National Institute for Nuclear Physics (INFN) and by the Department of Physics and Astronomy of the University of Catania. CSGF allows build, in a fast and easy way, standard-based web 2.0 Science Gateways to exponentially increase the number of potential users of Distributed Computing Infrastructure (DCIs) worldwide. The CSGF is released under the Apache 2.0 license and all code is available on GitHub².

The framework, conceived in the context of both EU-funded and national project, has been used in the last 5-6 years to develop Science Gateways for several EU co-funded projects such as: DECIDE³, EarthServer⁴, EUMEDGRID-Support⁵, GISELA⁶, DCH-RP⁷, INDICATE⁸ and CHAIN-REDS⁹, just to name a few. The framework is currently being completely re-engineered in the context of the INDIGO-DataCloud¹⁰ to offer a rich RESTful

¹ <http://www.catania-science-gateways.it/>

² <https://github.com/csgf/>

³ <https://www.eu-decide.eu/>

⁴ <http://www.earthserver.eu/>

⁵ <http://www.eumedgrid.eu/>

⁶ <http://www.gisela-grid.eu/>

⁷ <http://www.dch-rp.eu/>

⁸ <http://www.indicate-project.org/>

⁹ <https://www.chain-project.eu/>

¹⁰ <https://www.indigo-datacloud.eu/>

API as well as to include additional cloud, data and workflow data management functionalities and to improve its performances and reliability.

This short guide is intended for the users of the CSGF-based Science Gateway¹¹ built for the EGI Long Tail of Science¹² (LToS). Within these pages, users will find an introduction on how to use the scientific applications currently available in the portal. The guide does not cover the registration and authentication of users, for which the reader can see this page¹³.

Once the user is successfully registered, authenticated and authorized s(he) is then presented with the web page(s) of the application(s) s(he) is allowed to execute on the EGI LToS infrastructure.

The complete list of applications the user is entitled to run is available navigating the Applications menu located in the top bar, as shown in Figure:



2. Applications

In the next sub-sections, we will describe each of these applications in details.

2.1 R

The R Project for Statistical Computing¹⁴ is a language for statistical computing and graphics. It is a GNU project which is similar to the S language and environment which was developed at Bell Laboratories (formerly AT&T, now Lucent Technologies) by John Chambers and colleagues. R can be considered as a different implementation of S (see this page¹⁵ for more information).

R provides a wide variety of statistical (linear and nonlinear modelling, classical statistical tests, time-series analysis, classification, clustering, etc.) and graphical techniques, and it is highly extensible. The S language is often the

¹¹ <http://csgf.egi.eu>

¹² https://wiki.egi.eu/wiki/Long-tail_of_science

¹³ <https://access.egi.eu/start>

¹⁴ <https://www.r-project.org/>

¹⁵ <https://www.r-project.org/about.html>

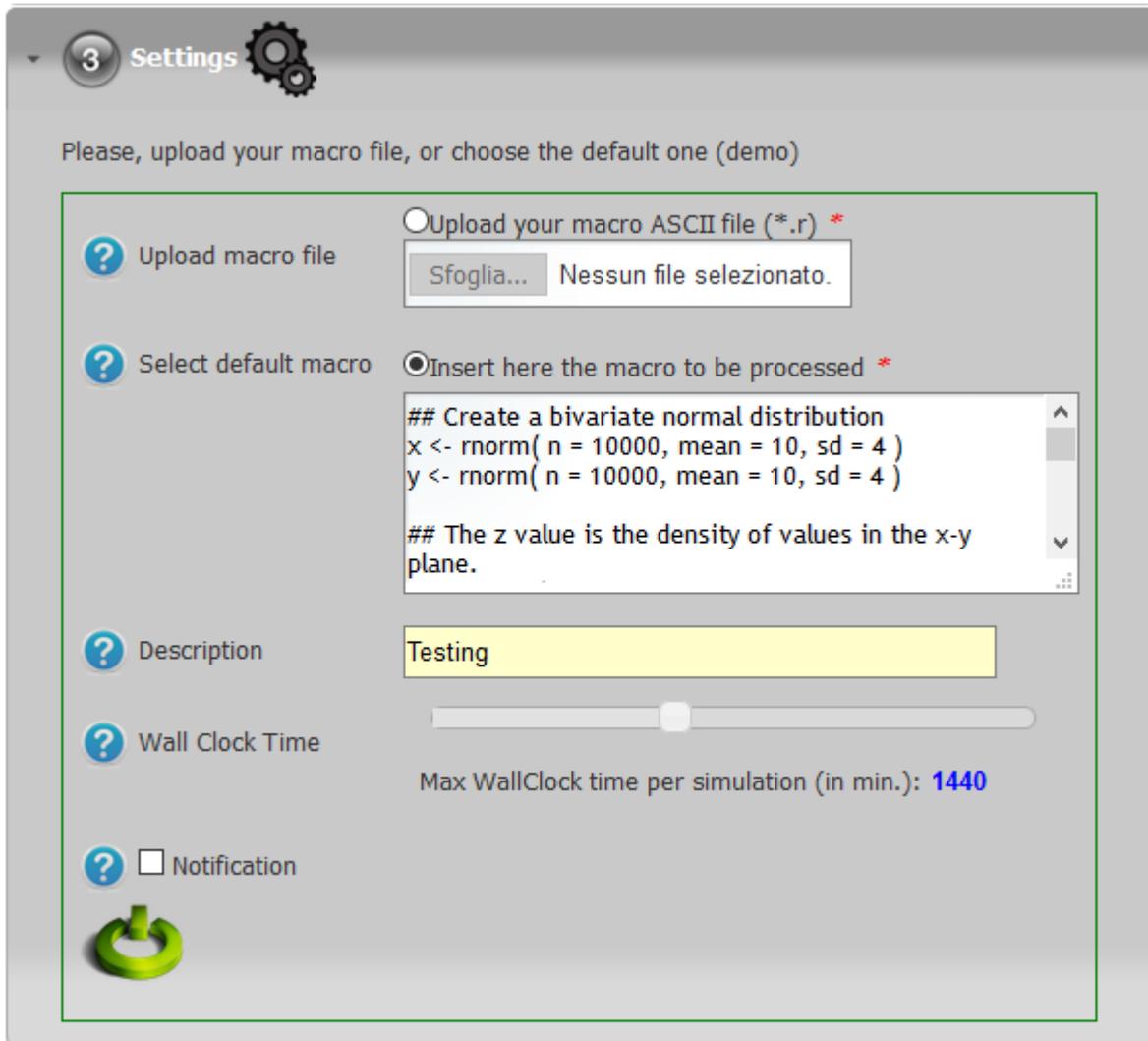
vehicle of choice for research in statistical methodology, and R provides an Open Source route to participation in that activity.

One of R's strengths is the ease with which well-designed publication-quality plots can be produced, including mathematical symbols and formulae where needed. Great care has been taken over the defaults for the minor design choices in graphics, but the user retains full control. R is available as Free Software under the terms of the Free Software Foundation's GNU General Public License in source code form. It compiles and runs on a wide variety of UNIX platforms and similar systems (including FreeBSD and Linux), Windows and MacOS.

A dedicated web page to access the service and execute the Statistical R for Computing is available at this URL¹⁶. The service is organized in three different accordions. The first one shows some generic information about the service, the second shows the list of computing resources entitled to run statistical analysis on High Throughput Computing (HTC) resources.

In the last accordion, the user can select to upload a local R macro file selecting the Browse button, or using the text box. By default, a macro demo is loaded when the corresponding radio button is clicked.

¹⁶ <https://csgf.egi.eu/run-r>



2.1 Chipster

Chipster is a user-friendly analysis software for high-throughput data. It contains over 300 analysis tools for next generation sequencing (NGS), microarray, proteomics and sequence data. Users can save and share automatic analysis workflows, and visualize data interactively using a built-in genome browser and many other visualizations. Chipster's client software uses Java Web Start to install itself automatically, and it connects to computing servers for the actual analysis. Please see the Chipster main site¹⁷ for courses, updates and other information. A dedicated web page to create a Chipster account for accessing the open source platform is available at this page¹⁸.

¹⁷ <http://chipster.csc.fi/>

¹⁸ <https://csgf.egi.eu/run-clustalw2>

2 Settings 

Please, specify your Chipster credential to access the open source platform

 Alias	larocca
 Password
 Re-type Password

 Get Notifications

 Password must meet the following requirements:

- At least **one letter**
- At least **one capital letter**
- At least **one number**
- At least **two special characters** (the character ':' is not permitted)
- Be at least **8 characters**

 Please activate notifications to get notified by e-mail



If 'Get Notification' is checked, the user will receive via e-mail instructions on how to access the Chipster Server.

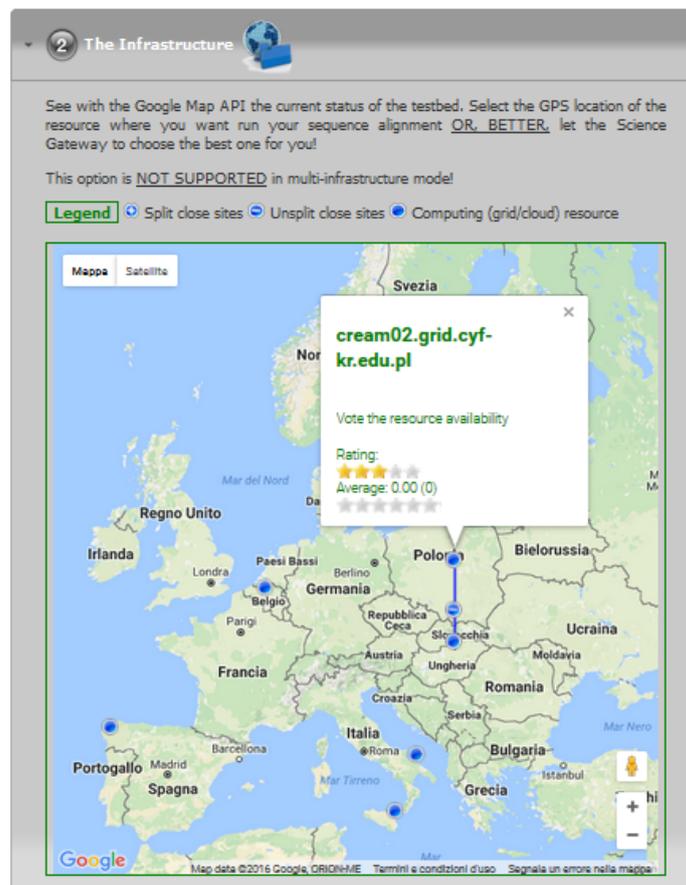
2.3 ClustalW2

ClustalW2¹⁹ is a widely used multiple program for multiple alignment of nucleic acid and protein sequences. sequence alignment computer program.

The program accepts a wide range on input formats including: NBRF/PIR, FASTA, EMBL/Swissprot, Clustal, GCC/MSF, GCG9 RSF, and GDE, and executes the following workflow:

- Pairwise alignment;
- Creation of a phylogenetic tree (or use a user-defined tree);
- Use of the phylogenetic tree to carry out a multiple alignment.

A dedicated web page to access the service and execute the Multi Sequence Alignments for DNA or protein is available at this page²⁰. The service is organized in three different accordions. The first one shows some generic information about the service, the second shows the list of computing resources entitled to run the multi sequence alignments process.



¹⁹ <http://www.clustal.org/>

²⁰ <https://csgf.egi.eu/run-clustalw2>

By default, the LToS CSG will send the request to one of these resources selected randomly, but the user can also select one of them if s(he) wants to force the execution of the sequence alignment on a specify one.

All the input settings needed to start the multi sequence alignment can be specified in the last accordion, as shown in the figure below.

The screenshot shows a web interface titled "Sequences & Options" with a gear icon. It contains several sections for configuring a sequence alignment job:

- Input Section:** A text area for the sequence to be analyzed, with a "Scogli file" button and a "Nessun file selezionato" message. The text area contains two sequences: `>seq0` and `>seq1`, each followed by a protein sequence. The first sequence is `FQTWEEFSRAAEKLYLADPMKVRVVLKYRHVDGNLCIKVTDDLVCVYRTDQAQDVKKIEKF` and the second is `KYRTWEEFTRAAEKLYQADPMKVRVVLKYRHCDGNLCIKVTDDVVCLLYRTDQAQDVKK`.
- Alignment Type:** Radio buttons for "Slow, but accurate" (selected) and "Fast, but approximate".
- Description:** A text input field with the placeholder "Please, insert here a description for your job".
- Notification:** A checkbox labeled "Notification" with an email icon.
- Step 2 - Slow Pairwise Alignment Options:** A section with a right-pointing arrow. It includes a note: "The default settings will fulfill the needs of most users. Configure advanced settings for the selected alignment options (if needed)." Below this are three dropdown menus: "Protein Weight Matrix" set to "Gonnet", "GAP OPEN" set to "10", and "GAP EXTENSION" set to "0.1".
- Step 3:** "Set 3 - Set up your Multiple Sequence Alignment Options" with a downward arrow.
- Step 4:** "Set 4 - Set up your Advanced Output Options" with a downward arrow.

At the bottom left, there is a green circular arrow icon.

Each simulation will produce the following files:

- `std.out`: the standard output file;
- `std.err`: the standard error file;
- `outputs.tar.gz`: the file containing the results of the Monte Carlo simulation.

2.4 Semantic Search Engine

The Semantic Search Engine (SSE) is a framework conceived to demonstrate the potential of information coupled with semantic web technologies to address the issues of data discovery and correlation. Two different version of the SSE are available on the CSG:

A simple Semantic Search Engine (SSE) customized for the DARIAH Competence Centre²¹ project which allows users to search in the e-Infrastructure Knowledge Base, in more than 100 languages across more than 30 million resources contained in the thousands of semantically enriched Open Access Document Repositories²² and Data Repositories²³. Search results are ranked according to the Ranking Web of Repositories²⁴.

The search can be performed:

- By typing one or more keywords in the textbox search. After entering the search string, clicking the “Search” button the system submit a SPARQL²⁵ query that retrieve the results from the e-Infrastructure Knowledge Base.

Example of query to retrieve the resources with the keyword in title.

```
SELECT distinct ?s WHERE {  
  
  ?s dc:title ?title. ?title bif:contains " + keyword + ". ?s  
<http://semanticweb.org/ontologies/2013/2/7/RepositoryOntology.owl#isResourceOf> ?rep. ?rep  
<http://www.semanticweb.org/ontologies/2013/2/7/RepositoryOntology.owl#rank>  
?rank. }ORDER BY ASC(?rank) limit 20 offset 0
```

- Using search filters. Inserting in textbox search string, for example, “author: G. Smith” and clicking “Search” will be launched a SPARQL query that will return only the resources of the e-Infrastructure Knowledge Base that have at least one author who matches “Smith G.”. The possible filters (according to the Dublin Core Standard) are : dc: author, dc: subject dc: type and dc: publisher. Clicking the “Examples”

²¹ https://wiki.ei.eu/wiki/Competence_centre_DARIAH

²² <http://www.sci-gaia.eu/e-infrastructures/knowledge-base/oadr-map/>

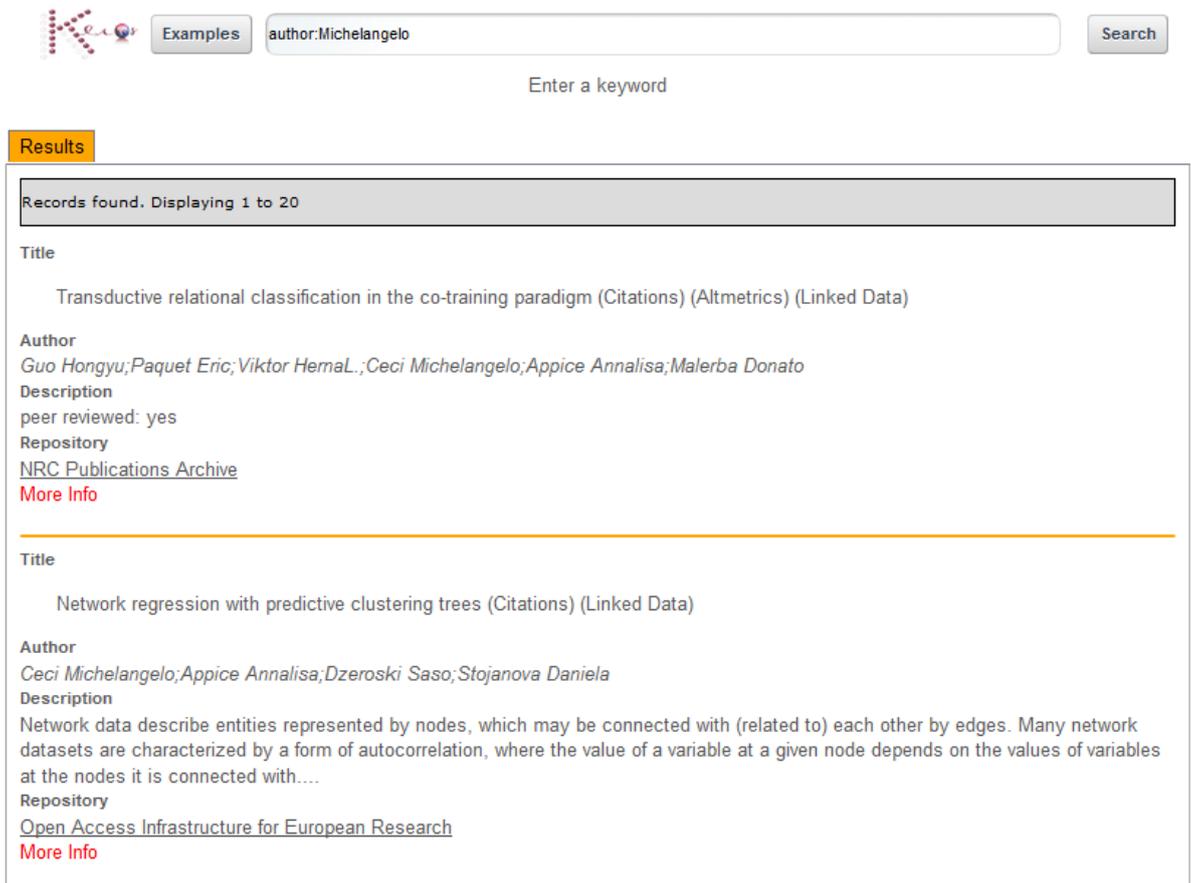
²³ <http://www.sci-gaia.eu/e-infrastructures/knowledge-base/data-repositories-map/>

²⁴ <http://repositories.webometrics.info/>

²⁵ <https://www.w3.org/TR/rdf-sparql-query/>

button are show examples of search filters in 4 of the 110 languages supported by the system.

In both cases, the results of the query will be processed further, to obtain the final results, as shown in the following figure.



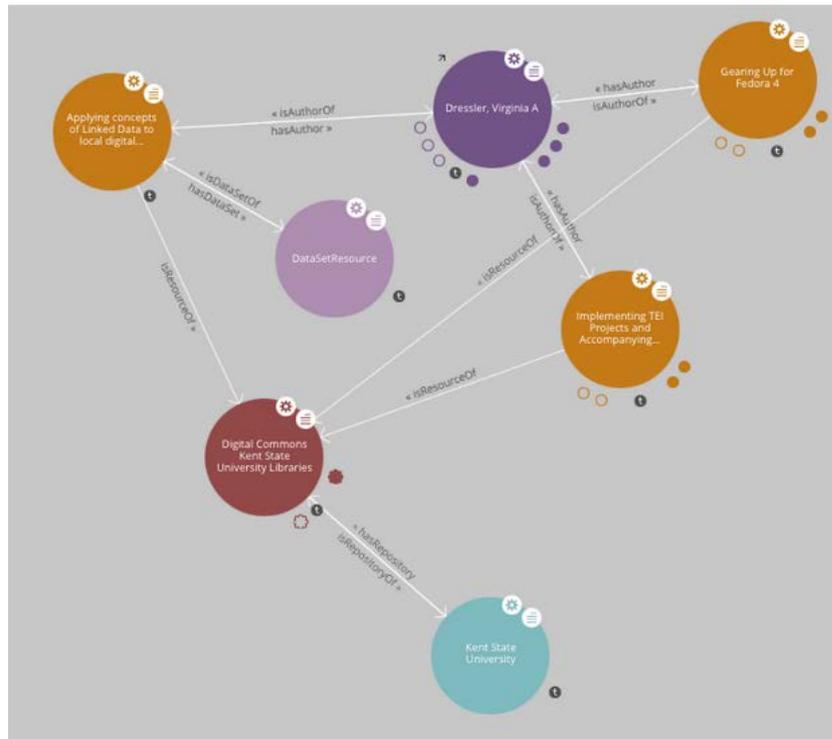
The screenshot shows a search interface with a search bar containing 'author:Michelangelo' and a 'Search' button. Below the search bar is a 'Results' section with a 'Records found. Displaying 1 to 20' header. The first result is titled 'Transductive relational classification in the co-training paradigm (Citations) (Altmetrics) (Linked Data)'. The author list includes Guo Hongyu, Paquet Eric, Viktor HemaL., Ceci Michelangelo, Appice Annalisa, and Malerba Donato. The description indicates it is peer reviewed and available in the NRC Publications Archive. A 'More Info' link is provided. The second result is titled 'Network regression with predictive clustering trees (Citations) (Linked Data)'. The author list includes Ceci Michelangelo, Appice Annalisa, Dzeroski Saso, and Stojanova Daniela. The description explains that network data describe entities represented by nodes connected by edges. The repository is the Open Access Infrastructure for European Research, and a 'More Info' link is also present.

Clicking on “[More Info](#)” link you can access all the details of the resource.

Clicking on “[Citations](#)” link, if available, are displayed more information from Google Scholar about the versions and the quotations of the resource.

Clicking on “[Linked Data](#)” link will be open a new tab that shows a LodLive²⁶-based graph of the resource with all its metadata.

²⁶ <http://en.lodlive.it/>



- A parallelised SSE (PSSE) configured to simultaneously search across the above e-Infrastructure Knowledge Base, Europeana²⁷, Cultura Italia²⁸, Isidore²⁹, OpenAgris³⁰, PubMed³¹ and DBpedia³² platforms (others can be added upon request to sg-licence@ct.infn.it).

²⁷ <http://www.europeana.eu/portal/>

²⁸ <http://dati.culturaitalia.it/?locale=it>

²⁹ <http://www.rechercheisidore.fr/>

³⁰ <http://agris.fao.org/openagris/index.do>

³¹ <http://www.ncbi.nlm.nih.gov/pubmed>

³² <http://wiki.dbpedia.org/>

The screenshot shows a search interface with a search bar containing 'bellis perennis' and a 'Search' button. Below the search bar are several database filters: E-INFRA-KB, OPENAGRIS, EUROPEANA, CULTURAITALIA, ISIDORE, PUBMED, DBPEDIA, and ENGAGE. The search results are displayed in a table with the following content:

Records found. Displaying 1 to 4

Title
 Prvo priopćenje o vrsti *Puccinia distincta* McAlpine, novoj europskoj hrđi na tratinčicama (*Bellis perennis* L.) iz Hrvatske (Citations) (Linked Data)
 First report of *Puccinia distincta* McAlpine, the new European rust on daisies (*Bellis perennis* L.), from Croatia (Citations) (Linked Data)

Author
 Weber Roland W. S.; Lehrbereich Biotechnologie Universität Kaiserslautern Paul-Ehrlich-Str. 23 67663 Kaiserslautern Germany; Jurc Dušan; Gozdarski inštitut Slovenije Večnapot 2 1000 Ljubljana p.p. 2985 Slovenia

Description
 Teška infekcija hrđe na divljim i uzgajanim tratinčicama (*Bellis perennis* L.) otkrivena je u travnju i svibnju 2000. na izoliranom priobalnom području između Lovrana i Opatije. Patogena gljiva bila je determinirana kao *Puccinia distincta* McAlpine, koja se protekle četiri godine brzo širi po Europi, ...

Repository
[Portal of scientific journals of Croatia](#)
[More Info](#)

3. MyJobs

To monitor active jobs and download results at the end of the computation, users just need click on the MyWorkSpace link which appears in the top bar.



A complete list of Active and Done jobs will be showed in a table view as shown in the Figure below. Statuses are automatically updated every 15 minutes so there is no need to reload this page more frequently.

Once your jobs have finished, user has 96 hours to retrieve the output. Beyond that time, the output of jobs will automatically be deleted from the Science Gateway in order not to fill its storage with undesired stuff.

Copy Print Save Download Job output Search:

Show 10 entries First Previous 1 Next Last

info job	Application Name	User Description	Started on (UTC)	Status
R		test5	2016-08-26 08:52:41.0	RUNNING
R		Simulation Started	2016-08-26 08:10:59.0	
R		Simulation Started	2016-08-29 14:08:52.0	SUBMITTED
R		Simulation Started	2016-08-29 14:08:18.0	SUBMITTED

Showing 1 to 4 of 4 entries First Previous 1 Next Last

4. Support

If you have any question, please contact the Catania Science Gateway support [team](#).