



- RO-LCG 2016, Bucharest-Magurele, 26-28 October 2016 -

Integration of HTC and HPC tools for solving complex problems in computational biology

**Dragos Ciobanu-Zabet¹, Dragos Hont², George Necula¹,
Ionut Vasile¹, Mihnea Dulea¹**

**¹ Dept. of Computational Physics & Information Technology
Horia Hulubei National Institute for R&D in Physics and Nuclear Engineering (IFIN-HH)**

² S.C. Totalsoft SA

Molecular modeling makes extensive use of:

- ❑ **HTC**, for database interrogations, massive data retrieval and conversion, high-throughput VLS, protein structure and function prediction, genomics, etc.
- ❑ **HPC**, for molecular dynamics simulations, chemical structure calculations, etc.

The molecular modeling of complex cellular subsystems implies procedures that require both sequential and parallel computing steps in order to obtain physically significant results, such as the description of drug-protein interaction in bacterial membranes.

The implementation of such a procedure often involves many different software tools for retrieval, conversion, processing and analysis of data, handled by experienced users.

We recently designed and currently implement a new integrated system for molecular modeling that automates the above-said procedures by means of workflows, in order to simplify the user's tasks.

The system consists of an extensible and scalable pool of **HTC and HPC resources** which are accessible to the user through a graphical frontend (**portal**).

The data handling and processing through various retrieval, conversion and modeling tools is managed by means of automatic, programmable and reusable workflows.

FUNCTIONALITY

PURPOSE: to empower the researchers with software tools allowing them to perform faster biological system setups, and numerical simulations.

The **portal** is currently customized for the modeling and simulation of subcellular structures of Gram-negative bacteria.

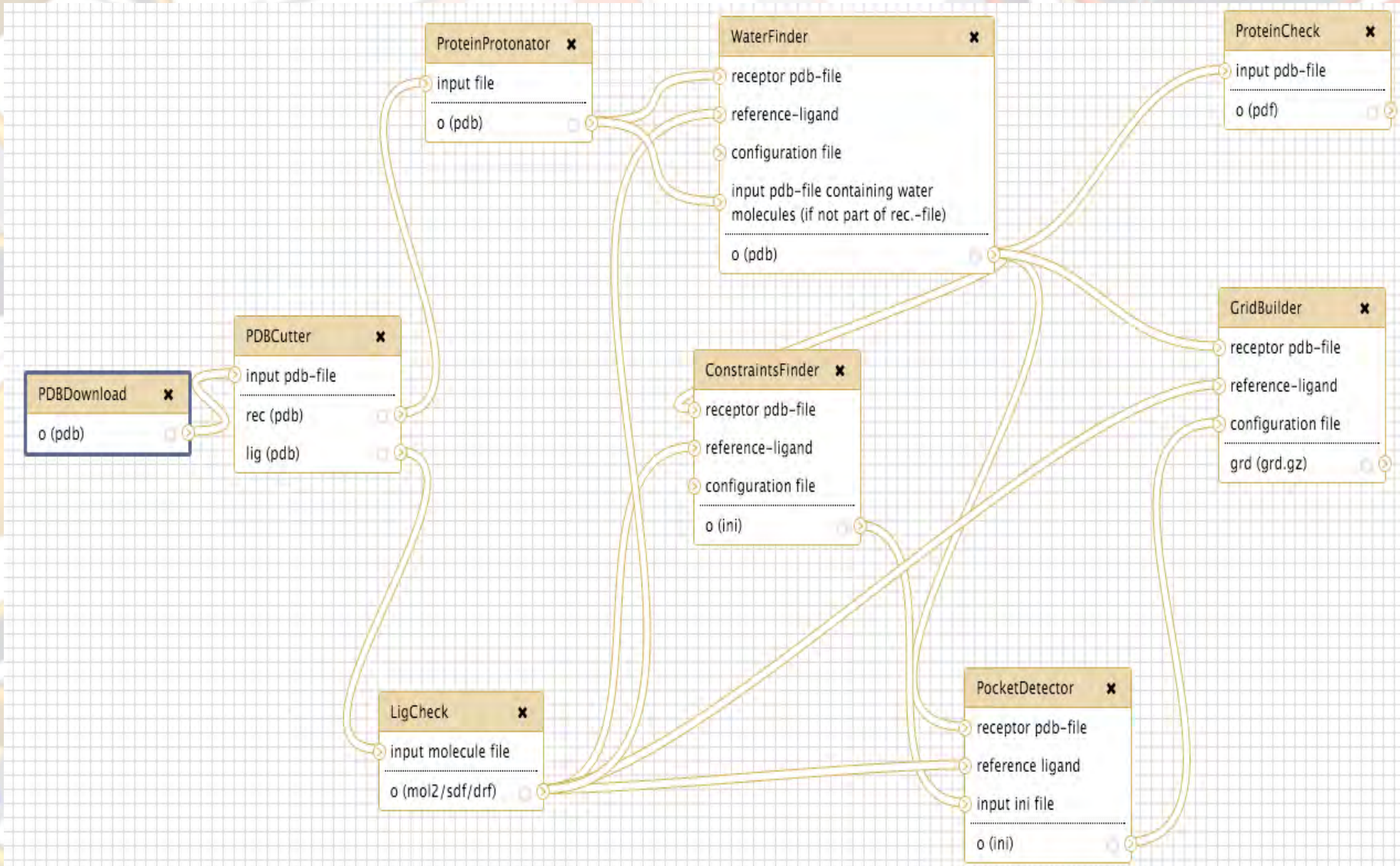
It provides access to software tools for microbiology and pharmacology researchers:

- computational modeling at molecular level
- workflows composed of tasks/modules that correspond to various modeling steps:

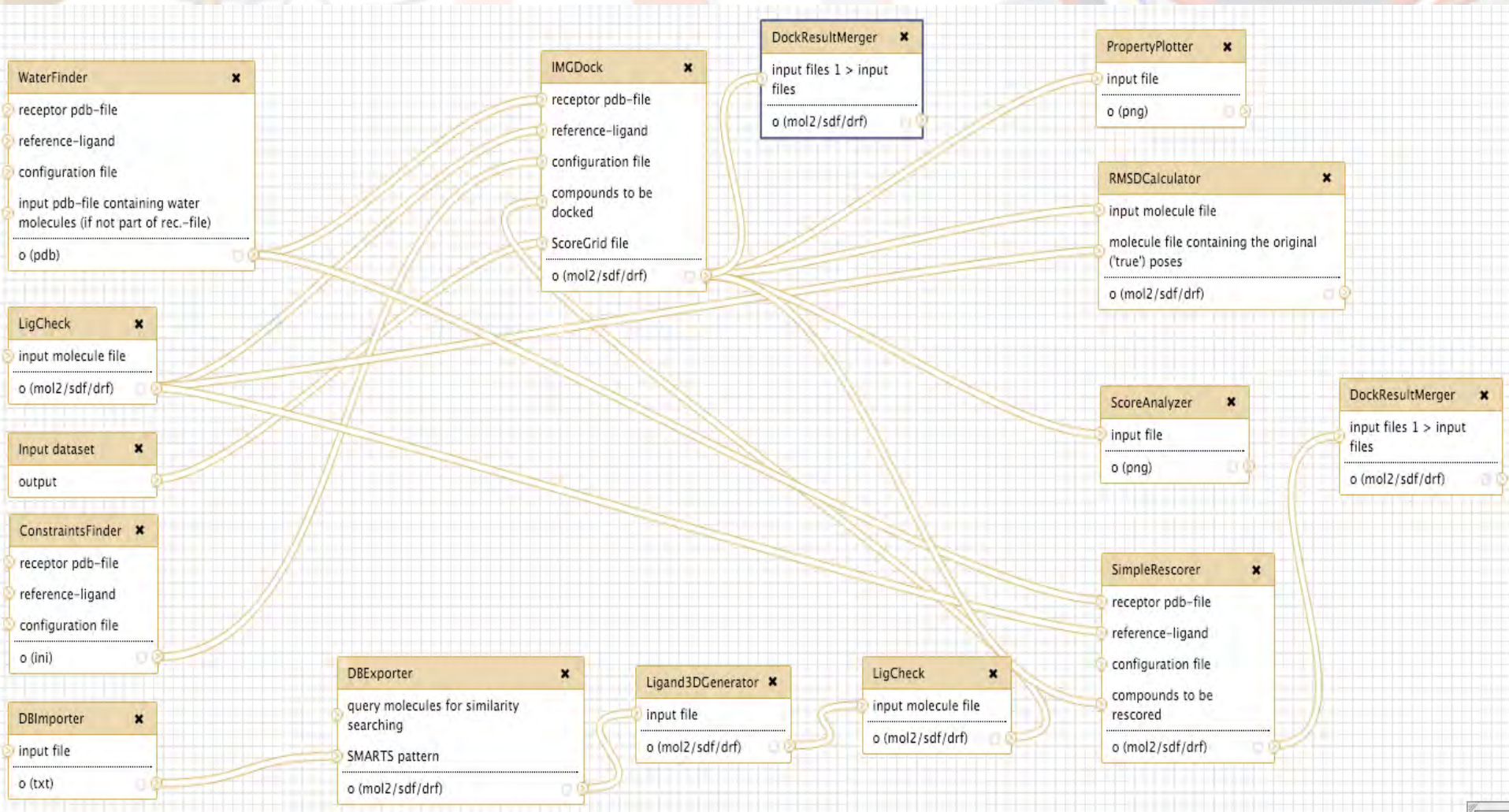
Examples:

- Search of molecular structures in databases (LipidBank, RCSB)
 - Data format conversion (OpenBabel)
 - 3D structures construction
 - Molecular modeling (CHARMM); molecular docking
 - Quantum chemistry (Gaussian, GAMESS)
- workflows are submitted for execution on the available HTC/HPC resources

Example1: Workflow for preparing data for molecular docking



Example2: Molecular docking workflow



DISTRIBUTED SYSTEM OF HTC & HPC RESOURCES

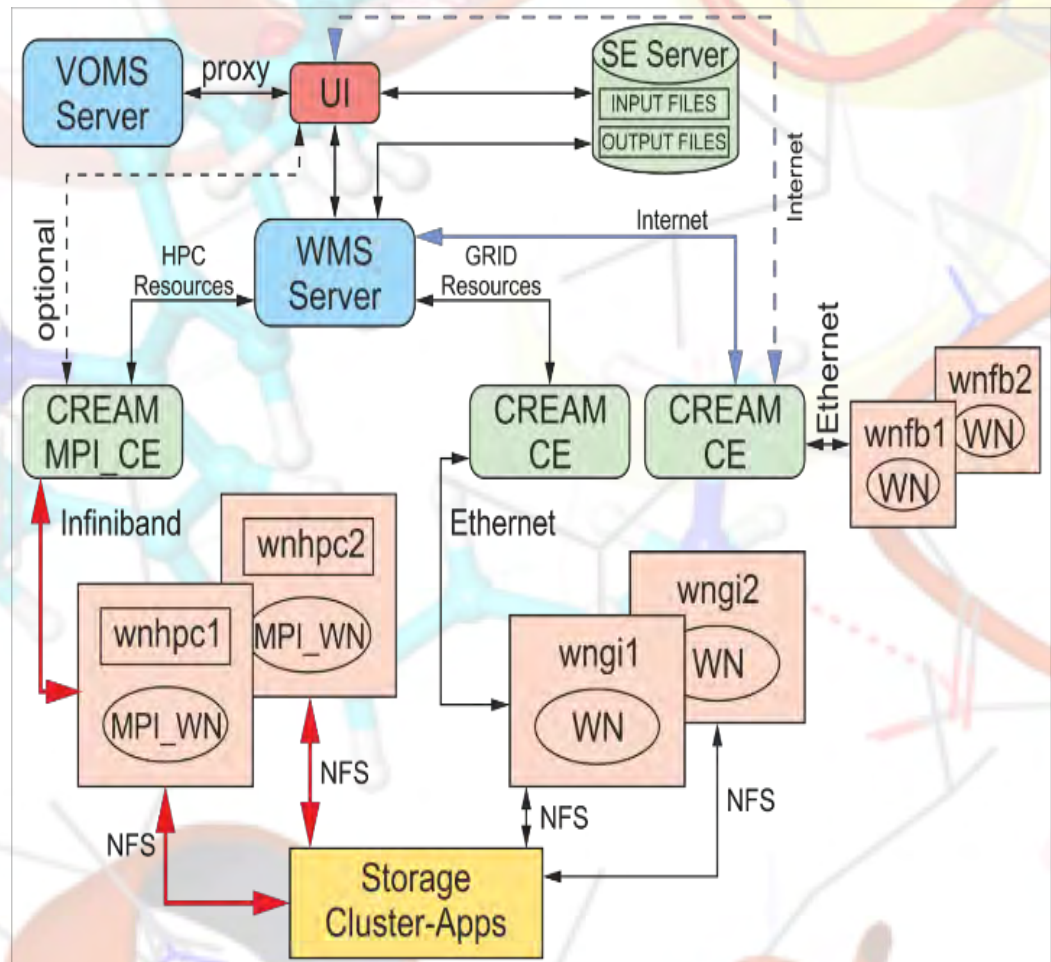
The system is currently using the resources of the GRIDIFIN site:

- HPC cluster
- HTC WNs
- Storage

GRIDIFIN: Two CEs (one for distributed and the other for parallel computing)

Multiple HTC/ HPC clusters can be attached to the system.

(The connection to another site is depicted in figure.)



IMPLEMENTATION

Grid jobs containing Taverna workflows are submitted and executed on the HTC/HPC clusters. The major advantage of Taverna is that it uses distributed services, optimizing the running of workflows, regardless of the local infrastructure.

The portal integrates graphical instruments for:

- Management of the input data
 - Download/upload of input files, jdl files, shell files, etc. from user or databases
 - Template files for creating jdl files, input for molecular dynamics, scripts
 - Edit and save input files
- Preparing, submitting for execution and management of jobs
 - Edit source files, jdl files, wrapper (for compiling), send jobs in execution
- Input data pre-processing and data analysis (MMTSB / AmberTools)
- Graphic representation and analysis of results (JSmol – in progress)

PORTAL: WORKFLOW MANAGEMENT

Example of simple workflow composed of two parts:

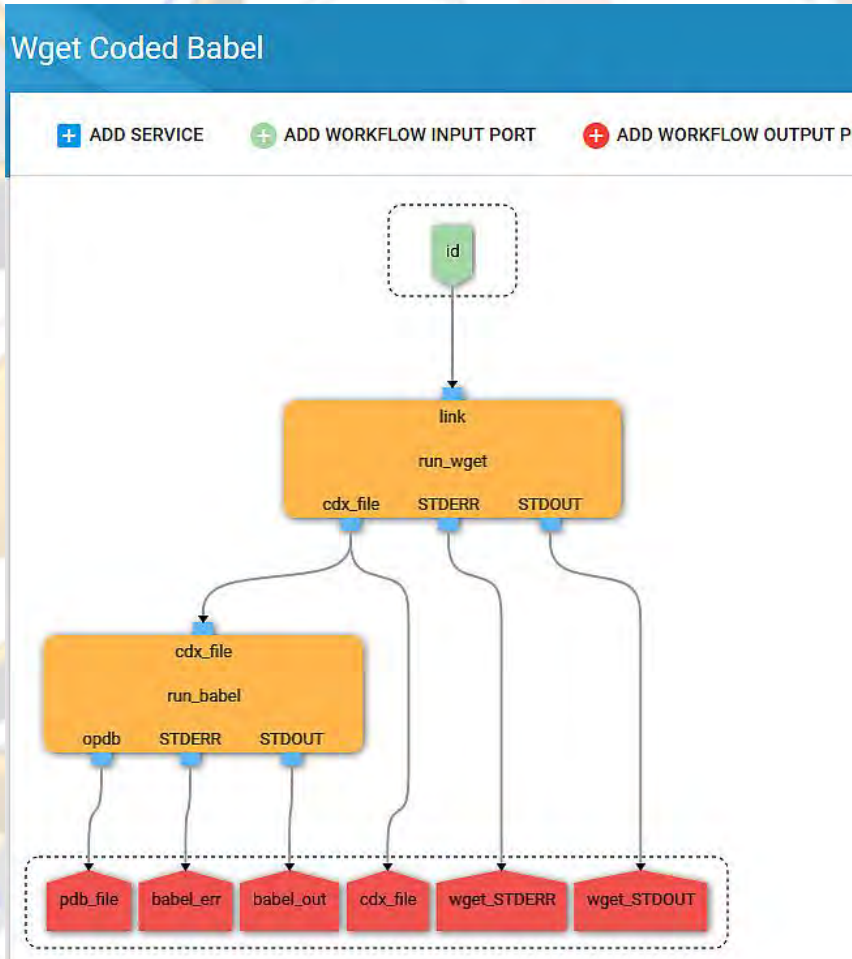
- file download from LipidBank
- file conversion from cdx format to pdb format using OpenBabel

The screenshot displays the SIMBAGRAN Workflows Management interface. On the left is a dark sidebar with navigation options: Home, Workflows (highlighted), Workflows to execute, Workflow execution history, and Users management. The main content area has a blue header with the title 'Workflows Management' and a search bar. Below the header, a red '+' button is visible. A dropdown menu is open, showing status options: All statuses, Under construction, Ready to execute, Ready to be publi..., In validation, and Rejected. The main area lists three workflow cards, each with a diagram and metadata:

- Workflow 1:** Date created: 10/06/2016, Created by: ion.stoenescu@fmail.com, Status: In validation.
- Workflow 2:** Wget Coded Babel, Date created: 09/26/2016, Created by: ion.stoenescu@fmail.com, Status: Ready to execute.
- Workflow 3:** Lipid Bank, Date created: 09/26/2016, Created by: ion.stoenescu@fmail.com, Status: Ready to execute.

WORKFLOW CREATION

Graphical representation and code



Wget Coded Babel

+ ADD SERVICE + ADD WORKFLOW INPUT PORT + ADD WORKFLOW OUTPUT

```
1 {
2   "dataflow": [
3     {
4       "name": "lipidbank",
5       "inputPorts": [
6         {
7           "annotations": {},
8           "granularDepth": "0",
9           "depth": "0",
10          "name": "id"
11        }
12      ],
13      "outputPorts": [
14        {
15          "annotations": {},
16          "lastPredictedDepth": "0",
17          "name": "wget_STDERR"
18        },
19        {
20          "annotations": {},
21          "lastPredictedDepth": "0",
22          "name": "wget STDOUT"
23        },
24        {
25          "annotations": {},
26          "lastPredictedDepth": "0",
27          "name": "cdx_file"
28        },
29        {
30          "annotations": {},
31          "lastPredictedDepth": "0",
32          "name": "pdb_file"
33        },
34        {
35          "annotations": {}
```

WORKFLOW EXECUTION HISTORY

The screenshot shows the SIMBAGRAN web interface for viewing workflow execution history. The left sidebar contains navigation options: Home, Workflows, Workflows to execute, Workflow execution history (selected), and Users management. The main content area is titled "Workflow execution history" and includes a search bar and filters for "All users", "All statuses", and "Workflow run date". Two workflow entries are displayed:

- wget coded babel initial**
 - Workflow run start date: 10/07/2016/11:50:52
 - Workflow run end date: 10/07/2016/11:56:56
 - Executed by: ion.stoenescu@mail.com
 - Arguments: id = G39_ideal.sdf; id_rcsb = 4c1w.pdb; center_x = 0; center_y = 0; center_z = 0; size_x = 80; size_y = 80; size_z = 80
 - Status: Executed with success
- Wget Coded Babel**
 - Workflow run start date: 09/30/2016/07:03:18
 - Workflow run end date: 09/30/2016/07:04:10
 - Executed by: ion.stoenescu@mail.com
 - Arguments: id = CLS0102.cdx
 - Status: Executed with success

Summary and conclusions

Graphical frontend (portal)

- Makes use of advanced (HTC/HPC) computing resources
- Zero user knowledge of HTC/HPC management
- Based on Taverna workflow management system
- Many tools in a single place
- Ease of use
- Graphical simplicity (user friendly)

THANK YOU FOR YOUR ATTENTION !