

“PORTAL access to Grid”

A proposal to TCG

Christophe Blanchet

CNRS IBCP

Chair of EGEE Bioinformatics activity

PORTALwg Kick-off

May, 2nd-3rd, 2007, CERN

- **PORTAL working group**
 - chair: C. Blanchet
- **Goals:**
 - The main objective of this "PORTAL" working group is to propose "best-practice" rules for the access of portals to the grid. By access to the grid, we mean that the portal should be able to store data and run job on the grid by delegation of real users, or by itself with its own credentials. To do so, a portal responsible should be able to get a server certificate for his portal from the EUGridPMA CAs, and then to be able to register this portal certificate to a VO allowed on the grid. Once the portal have been accepted into the concerned VO, it should be able to store and access data inside the VO area, and also to run job on site accepting this VO. The usage of these grid resources by the portal can be done by delegation of a grid- and VO-registered user, or as a resource provider for an user not registered in the grid.

- **Having agreed**
 - Christophe Blanchet , Bioinformatics usecase and chair
 - Hugues Benoit-Catin , Medical imaging usecase
 - Massimo Lamana , repr. of HEP experiments (repr. tbd)
 - Gidon Moont , GridPP portal usecase
 - Cal Loomis , small site repr. and representative of TCG
 - Rolf Rumler , large site repr. (repr. tbd)
 - Peter Kunszt , large site repr.
 - Dave Kelsey , repr. of JPSG
 - David Groep , repr. of EUGridPMA
 - Ake Edlund , repr. of MWSG (repr. tbd)
 - Jens Jensen , expert about authN/authZ
 - Christoph Witzig , expert about Shibboleth in SWITCH
- **Others ?**
 - Generic portal: No answer neither from Pgrade nor Genius
 - Jim Basney (US)

- **Interest ?**
 - PORTALwg has been announced
 - EGEE Bioinformatics meeting 3 (Valencia, March 6th, 2007)
 - § great interest of the Bioinformatics community
 - EGEE JPSG meeting, March, 13rd-14th, Cern
 - § good interest on the site operation and the policy aspects
 - other communities: Geophysic, Finance, ...

- **Collaboration with other projects**
 - OSG : certainly, according to last MWSG (March 1st-2nd), Jim Basney to be contacted
 - SwissGrid: participation of CSCS and Switch

- **Timeline**
 - from May 2007 to May 2008 ?

**PORTALwg
Kick-Off Meeting
May, 2nd-3rd, 2007
Cern**

- **Detail participant contribution**
- **Finalize the mandate of the PORTAL group**
- **Define a timeline to propose to TCG**
- **Suggest other participants if required.**
- **Work :-)**

- **Questions**

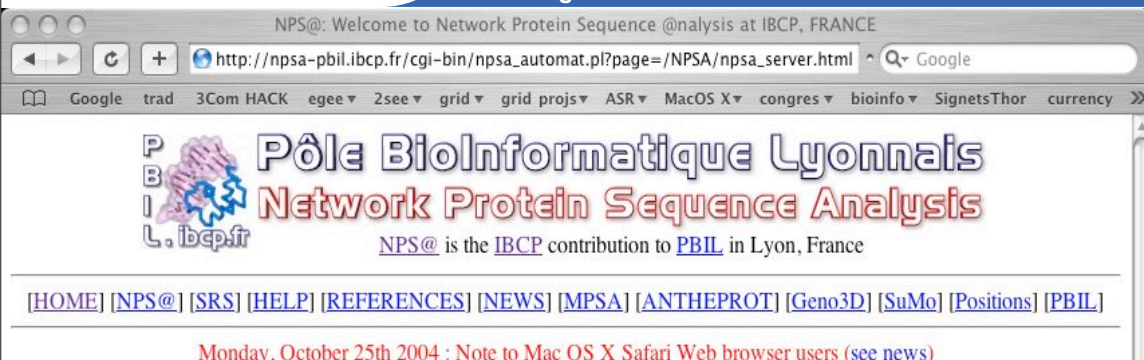
- How does one obtain service certificates from the EUGridPMA CAs?
- How to register those certificates in a VO?
- How to handle proxies of this certificate in a portal (e.g. wanting to add VOMS groups)?
- Which software should be allow to be used ?
 - § enabling users to put his own software
 - § restrict to only portal ones
- Do sites need to be advised of the use of these certificates (if so, how to inform them)?
- What are the accounting/logging practices to follow on the portal?
- Does the use of these need to appear in the VO AUP or similar documents?

- **Other questions ?**

- **May, 2nd, 14h00-17h00**
- **Introduction (C. Blanchet)**
- **Application usecases**
 - Christophe Blanchet "Bioinformatics GPSA portal"
 - Gidon Moont "GridPP portal"
- **Grid and Site operation**
 - local site and global platform constraints and security flaws
 - Cal Loomis "Grid site: CNRS LAL"
- **Discussion**

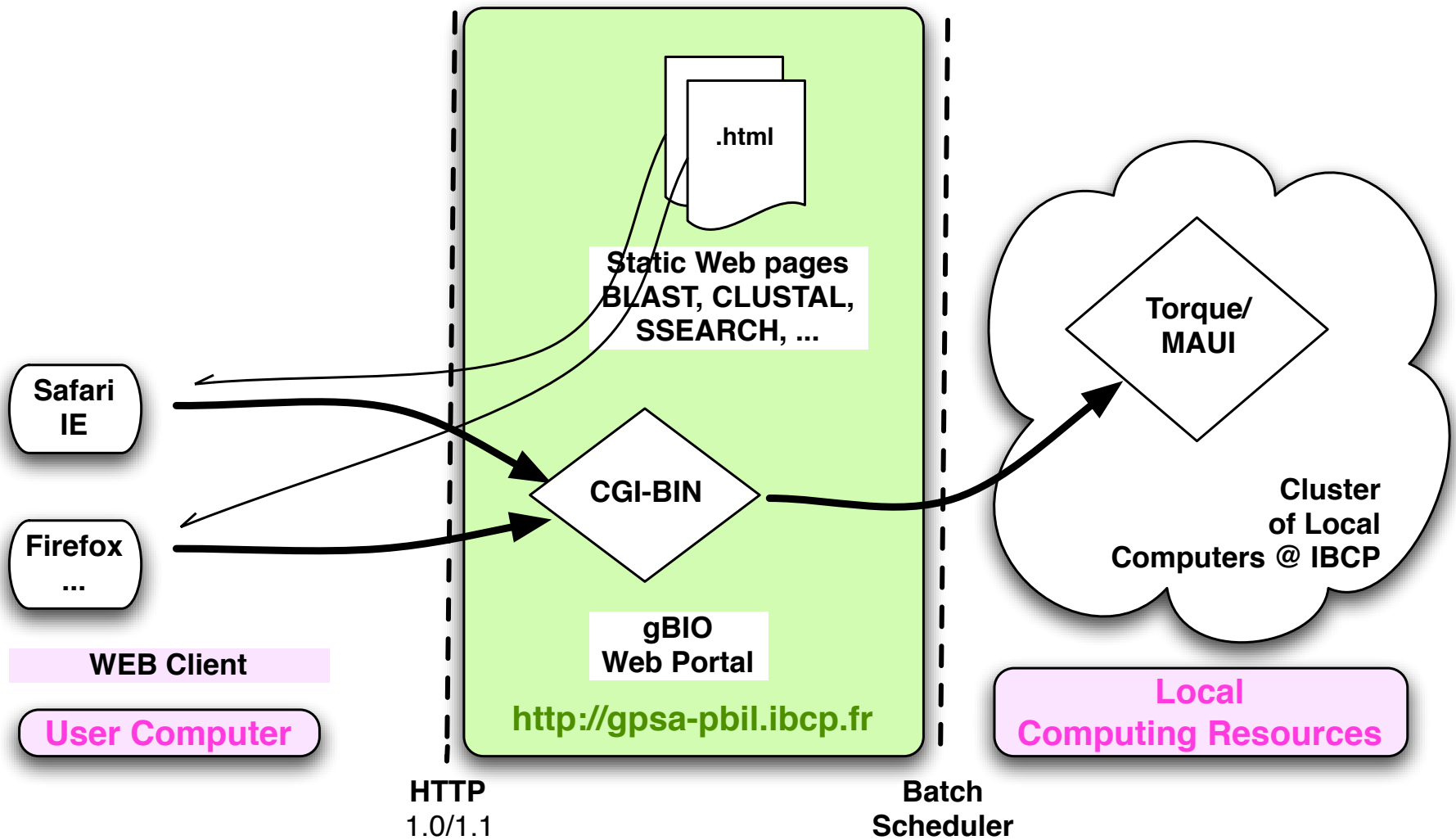
- **May, 3rd, 9h30-12h00**
- **Authentication and authorization: portal certificate, certificate authority and credential delegation**
 - David Groep "EUGridPMA: certificate authority at the European scale"
 - Jens Jensen "Single Sign-On and Identity Management on Grid"
- **Best practice and policy agreement:**
 - David Kelsey "JSPG Best practice and policy agreement on EGEE "
- **Mandate and timeline agreement**
- **Discussion and AOB**

UseCase-1 : Bioinformatics Web Portal



- [What is NPS@ ?](#)
- [Software facilities](#) to analyse NPS@'s data: [AnTheProt](#) and [MPSA](#).
- [Work with your own database](#)
- [Geno3D : Automatic modeling of proteins 3D structure](#)
- [SRS : Sequence Retrieval System](#)
- [Sequence homology search against proteic databases :](#)
 - [BLAST search](#) (protein (blastp) or nucleic (blastx) query sequence)
 - [PSI-BLAST search](#) (protein query sequence)
 - [FASTA search](#) (protein query sequence)
 - [SSEARCH search](#) (protein query sequence)
 - [HMMSEARCH](#) (protein query profile, hmmer format) **NEW**
- [Patterns and signatures search :](#)
 - [PATTINPROT](#): scan a protein sequence or a protein database for one or several p
 - [PROSCAN](#): scan a sequence for sites/signatures against PROSITE database
 - [InterProScan](#): scan a sequence for signatures against InterPro database
- [Profile building :](#)
 - [HMMBUILD](#): build a profile with HMMER (HMMER profile format) **NEW**
- [Multiple alignment:](#)
 - [Clustal W Protein](#) sequences (Des Higgins, EBI, Hinxton Hall, UK)
 - [Clustal W DNA](#) sequences (Des Higgins, EBI, Hinxton Hall, UK)

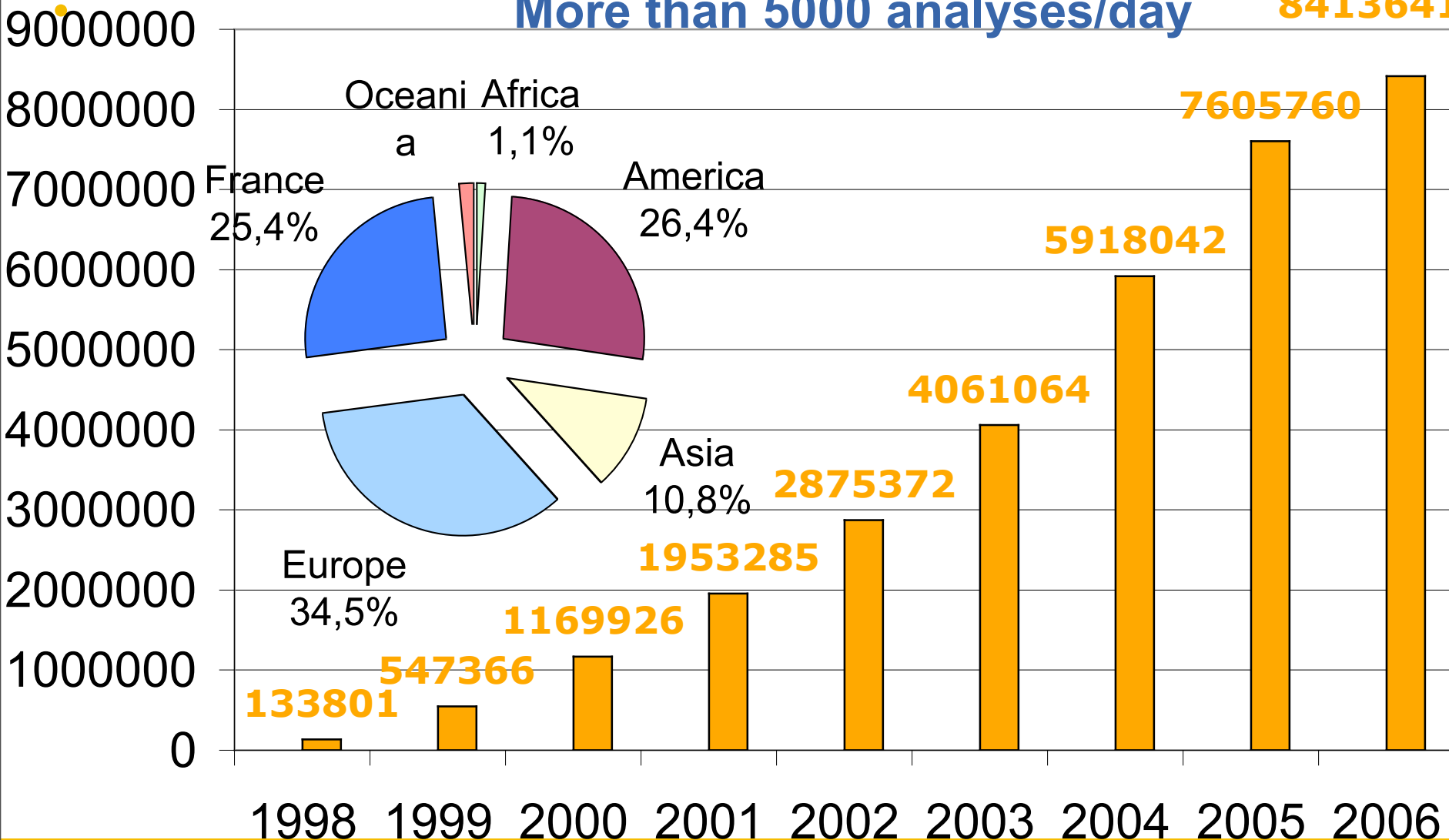
- **Network Protein Sequence Analysis (NPS@ release 3)**
<http://npsa-pbil.ibcp.fr>
- **Online since 1998**
- **46 integrated methods for protein sequence analysis**
- **12 Online up-to-date biological databanks**
- **International pointers: Expasy (Ch) , University of California, ...**
- **Ref.: “ NPS@: Network Protein Sequence Analysis”, Combet C., Blanchet C., Geourjon C. et Deléage G. Tibs, 2000, 25, 147-150.**



More than 8 millions analyses since 1998

More than 5000 analyses/day

8413641



Scientific objectives

- Molecular Bioinformatics: protein sequence analysis
- Analyse data from high-throughput Biology: complete genome projects, EST, complete proteomes, structural biology,
- Integration of biological data and tools

Method

- Provide Biologists with an usual Web interface: NPS@
 - § NPS@ Web portal online since 1998
 - § 46 tools & 12 updated databases
 - § + 9,000,000 jobs & 5,000 jobs/day
- Ease the access to updated databases and algorithms.
- Protein databases are stored on grid storage as flat files.
- Legacy bioinformatics applications
 - § Wrapping usual binary in grid environment
 - § transparent remote access with local filesystem
- Display results in graphical Web interface.

Status: Prototype




Institut de Biologie et Chimie des Protéines

Grid-enabled bioinformatics resources

- 9 algorithms
- 3 protein databases

Bioinformatics descriptors

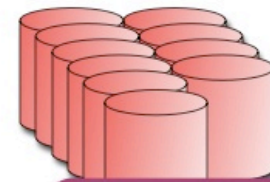
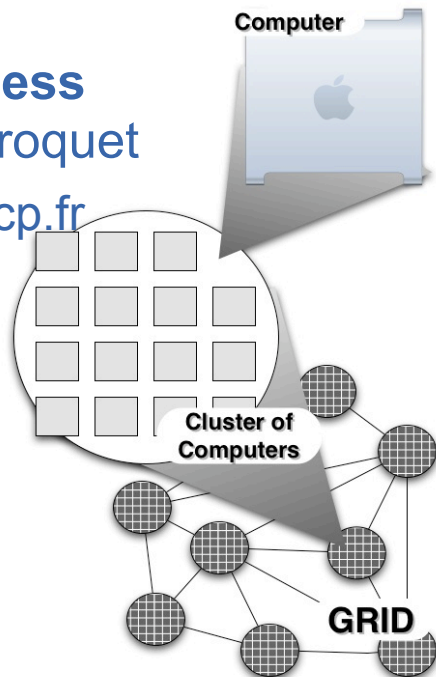
- XML framework, WSRF

Encryption system: EncFile

- Security: AES, Key sharing, M-of-N

Transparent access to grid files: Perroquet

<http://gpsa-pbil.ibcp.fr>



Databases



Software

Grid-enabling Bioinformatics

Legacy Bioinformatics Applications :

- * Wrapping tools with XML descriptors
- * BLAST, SSEARCH, FASTA, ClustalW, MultAlin, PattInProt, ...

Distributed databases:

- * **Encrypting Data** with EncFile system

* Swiss-Prot:

lfn:/grid/biomed/db/swissprot/last/sprot.fas

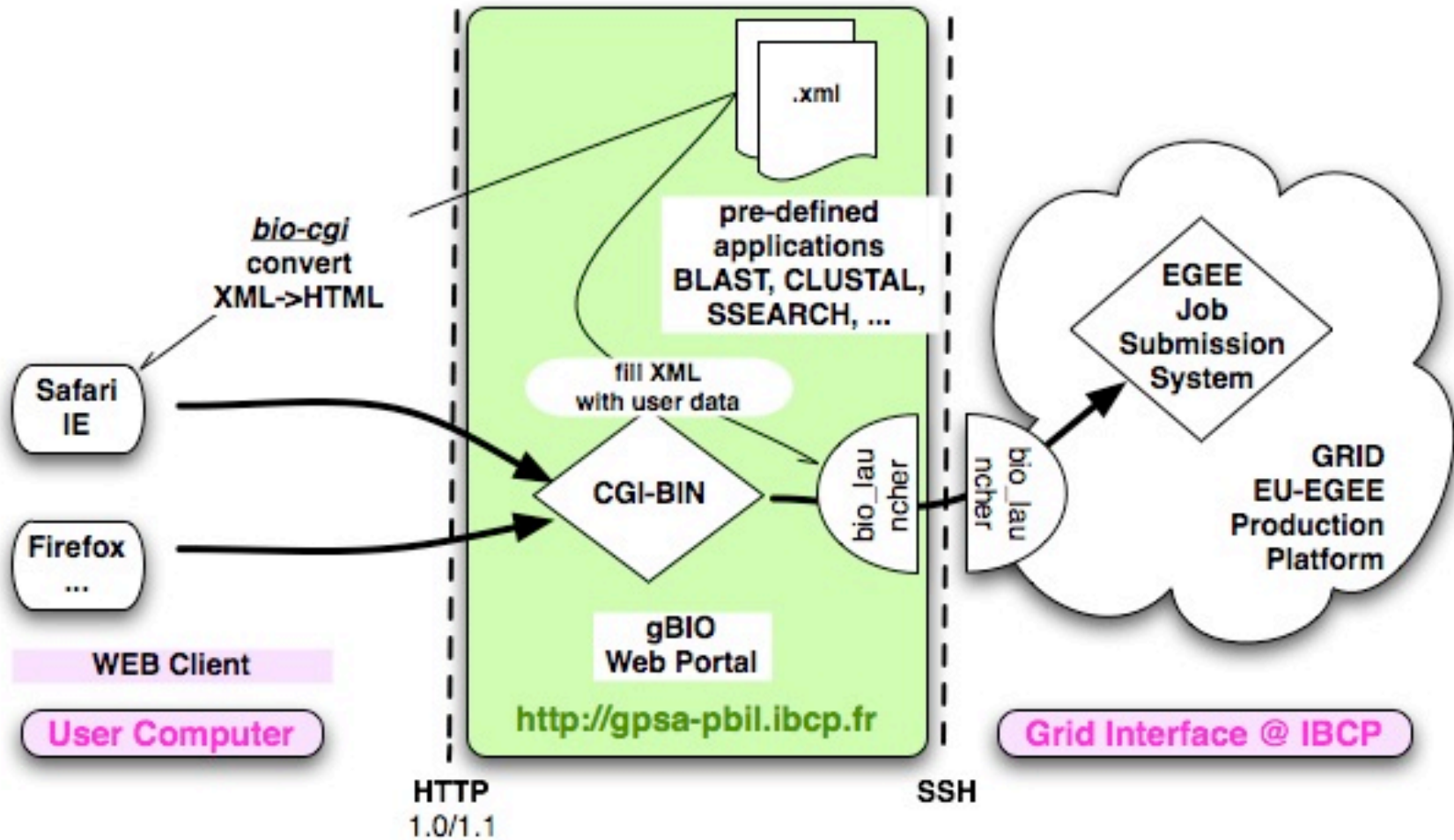
Symlink to (...)swissprot/50/4/sprot.fas

* TrEMBL:

lfn:/grid/biomed/db/trembl/33/4/spr.fas

Details : <http://gbio.ibcp.fr>





UseCase-2 : Bioinformatics Web Services

Advanced model explorer

Workflow Metadata for 'gBIOclustalwGrid'

Load Save New subworkflow Offline Reset

Workflow object	Retrie	Delay	Backof	Thread	Critica
Workflow model					
Workflow inputs					
sequences					
Workflow outputs					
alignment					
Processors					
gBIOclustalwGrid	0	0	1	1	
sequence-bank 'text/plain'					
attachmentList l("")					
result 'text/plain'					
Data links					
sequences-gBIOclustalwGrid:sequence-bank					
gBIOclustalwGrid:result-alignment					
Control links					

Available services

Search Watch loads

- Available Processors
 - Local Services
 - Biomart service @ http://www.biomart.org/
 - WSDL @ http://gbio.ibcp.fr/ws/gBIO.wsdl
 - porttype: gBioWSPortType [RPC]
 - gBIOclustalw
 - gBIOclustalwGrid**

Enactor invocation

Save as XML Save to disk Save to disk as website Excel

Status Results Process report

alignment

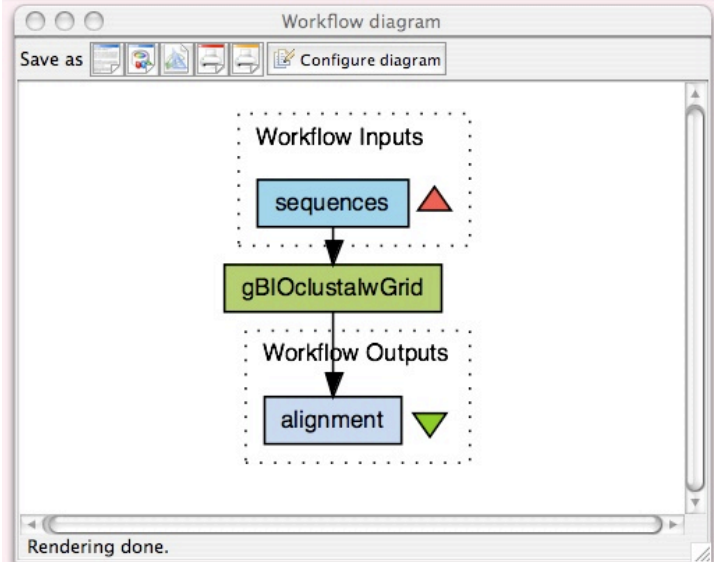
```

CLUSTAL W (1.83) multiple sequence alignment

CCPA_STRM5      MNTDDTITYDVAREAGVSMATVSRVVGNGK---NVKENTRKKVLEVIDRLD
DEGA_BACSU      ----MKTTIYDVAKAAGVSIITTVSRVINNTG---RISDKTRQKVMVMNEMA
FRUR_ECOLI      -----MKLDEIARLAGVSRRTTASYVINGKAKQYRVSDKTVEKVMVREHNY
RBTR_KLEAE      ---MKKITIYDLAELSGVSASAVSAILNGNWKRRISAKLAEKVTRIAEQG
                :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :
                :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :

CCPA_STRM5      GLASKKTTTVGVVPIPNIANAYFSLAKGIDDIATMYKYNIVLASSDEDDDK
DEGA_BACSU      ALTKRRTNMIALVAPDISNPFYGEIASEQIPIAMISQDKPLLPMDIVVIDV
FRUR_ECOLI      GLRAGRTRISGLVIPDLENTSYTRIANYLERQARQYQLLIACSEDQPDNE
RBTR_KLEAE      MLRSKSHVIGMIIPKYDNRVFGSIAERFEEARERGLLPITICTRRRPELE
                *  *  *  *  *  *  *  *  *  *  *  *  *  *  *  *  *  *  *  *  *  *

CCPA_STRM5      AKQVDGIIFMG---HHLTEKIRAEFSRARTPVVLSGTVDLHEQLPFSVNI
DEGA_BACSU      QKKVDGIIFATGIESHDSMSALEEIASAQIPIAMISQDKPLLPMDIVVIDV
FRUR_ECOLI      QRQVDALIVSTS---LPPHEFPYQRWANDPFPIVALDRALDREHFTSVVGADQ
RBTR_KLEAE      SWQVDVWVATG---ATNPDKISALCQQAGVPTVNLDPG---SLSPSVISDN
                **  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :  :
    
```



Load Inputs New Input

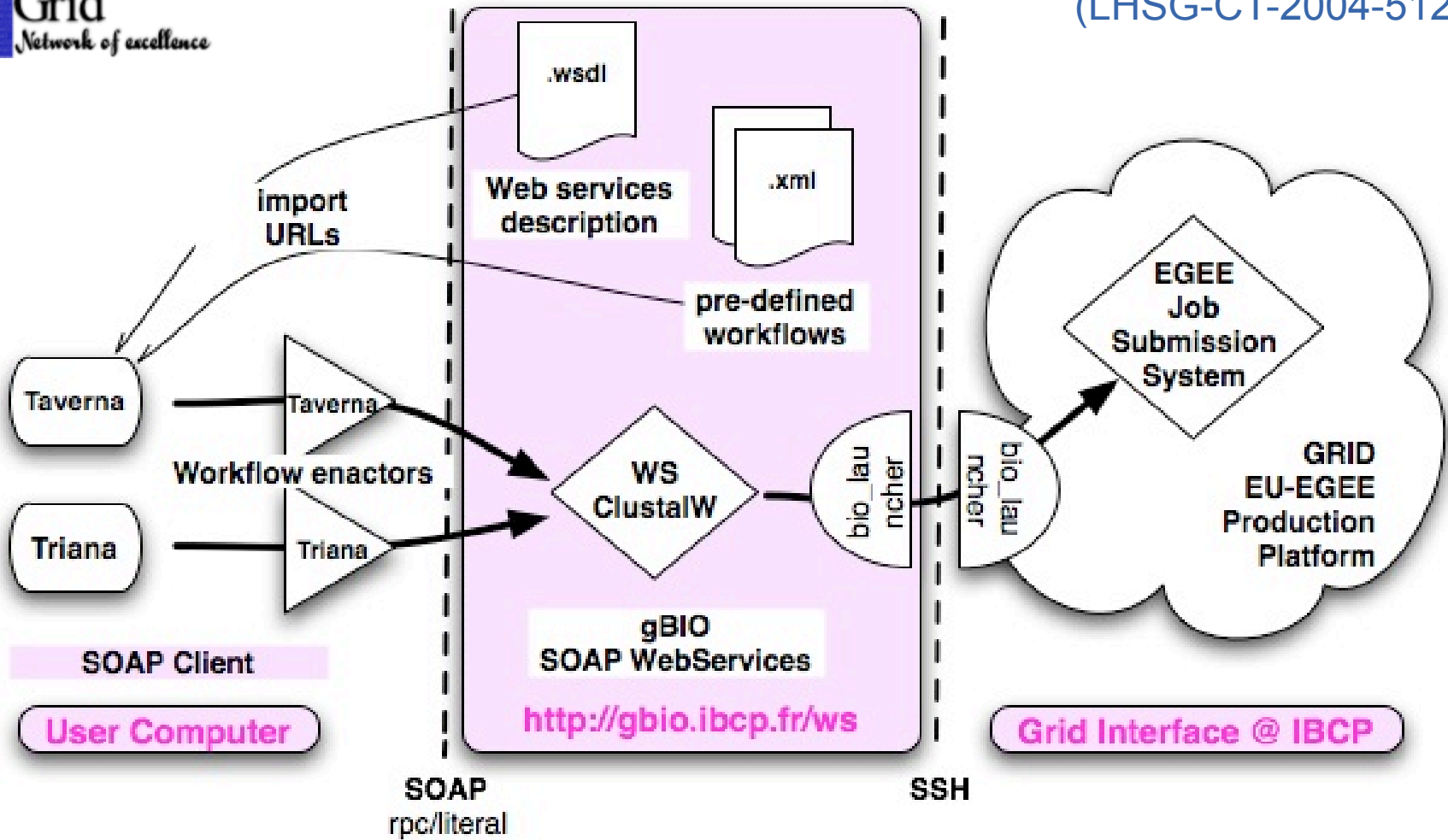
Input Document sequences Click to e

```

VDGIIFMGI
DKIAFVSGI
ATAAYVAE
VAMRMLTI
>DEGA_BA
MKTTIYDV
NMIALVAPDISNPFYGEIASEQIPIAMISQDKPLLPMDIVVIDVDRGGYEA
IFATGIESHDSMSALEEIASAQIPIAMISQDKPLLPMDIVVIDVDRGGYEA
TNIACIIGDSTTGTEKNRIKGRQAMEEAGVPIDESLIQTRFSLGKEEAGKLLDR
PTAIFAFNDVLACAAIQAARIRKIVPDDLIIIGFDNTILAEMAAPLTTVAQPIKE
ERHRTAGRSNRGRKAKQKIVLPELVVRHSTPLNT
>FRUR_ECOLI FRUCTOSE REPRESSOR.
MKLDEIARLAGVSRRTTASYVINGKAKQYRVSDKTVEKVMVREHNYHPNAVAA
TRISGLVIPDLENTSYTRIANYLERQARQYQLLIACSEDQPDNEMRCIEHLLQR
IIVSTLPPHEFPYQRWANDPFPIVALDRALDREHFTSVVGADQDDAEMLAELRL
TVLYLGAIPVSEI REOGFRTAWKDDPREVHEL YANSYERFAAAOIFKWIETI
    
```

Run Workflow

Call gBIO Web Services from Taverna, Triana, and other workflow tools



*Christophe Blanchet, Christophe Combet, Vladimir Daric and Gilbert Deléage
Web Services Interface to run Protein Sequence Tools on Grid, Testcase of Protein Sequence Alignment.
Lecture Notes in Computer Science : Biological and Medical Data Analysis (2006) 4345, 240-9*