

# Setting up a Bioinformatics Service Centre in a distributed environment

Patricia Rodriguez-Tomé

CRS4 Bioinformatica, Pula, Italy

<http://www.bioinformatica.crs4.org>

---

---

*Where are we ?*



# *Polaris*

## ❖ ICT:

- ✓ CRS4, Centro Sviluppo Materiali spa, UNICA – I3Lab, ICT Farm...

## ❖ BioMed:

- ✓ Neuroscienze PharmaNess scarl, Shardna spa, Biofarmitalia spa, Bio-Ker srl, Bioincubatore...

→ This is a good environment for a Bioinformatics group

→ 1 person from September 05 to mid February 06, 12 today...

---

---

# *Our users*

- ❖ Polaris users are our beta test users because of their proximity
  - ✓ Have small bioinformatics teams
    - We will provide advanced support
  - ✓ Have no team and/or low computing knowledge
    - We will provide basic and advanced support
- ❖ We have begun joint collaborations on new software development and data analysis

# *We work with/for*

- ❖ Many research projects/groups, bioscientists
    - ✓ Internal CRS4 Bioinformatics research group
      - Proteins, genome analysis, comparative genomics
    - ✓ Genotyping, Micro Arrays
      - Genetics, gene expression
    - ✓ Proteomics
      - Gel analysis, mass spectrometry
    - ✓ Data & text mining
  - ❖ Many questions, many databases, many programs
- 
-

# *A new genotyping lab*

- ❖ Medium throughput capacity data production
    - ✓ Databases: logbook (LIMS), results
    - ✓ Pipelines and Data Flow
  - ❖ Data Management
  - ❖ Data exploration
    - ✓ Experimental Data Analysis tools
    - ✓ Statistical Analysis for DNA chips
    - ✓ Workflows ?
  - ❖ Data security, privacy
- 
-

# *What we (will) provide*

- ❖ FTP server
  - ✓ Mirror of EBI, swissprot@exapsy, nr@NCBI, Ensembl
  - ✓ Our own (future) developments
- ❖ Database access: GO (FueGo), Ensembl
- ❖ Web server (BioPortale) gives compute access to :  
Fasta/Blast, Clustalw/Muscle/T-Coffee,  
genotyping, micro array and proteomics tools

# *Our setup*

- ❖ Independent entities with their own private network.
  - ❖ CRS4 has computing power. For Bioinformatics we can use:
    - ✓ A cluster of 24 nodes (dual AMD opteron) as a file server (20TB)
    - ✓ A cluster of 48 nodes (dual AMD opteron) as compute server
    - ✓ Other file servers, web /FTP servers, our desktops...
- 
-



# *No legacy*

- ❖ We are at the beginning of the project, we have no legacy
    - ✓ good isn't it ?
    - ✓ what can we use that will take us into the future ?
      - and prevent legacy for the next few years, hopefully
  - ❖ We need too use the machines available to us
  - ❖ “buzz words” we hear / read / don't always (fully) understand ...
    - ✓ Distributed computing
    - ✓ Grid
    - ✓ Workflows
- 
-

# *Distributed computing*

- ❖ It looks like us:
    - ✓ Distributed groups
    - ✓ Distributed computers
  - ❖ What goes into distributed computing
    - ✓ Batch processing, high performance computing
    - ✓ “the GRID”
    - ✓ PBS, Globus, LSF, IBM's Data Grid, Sun Grid Engine
    - ....
  - ❖ Distributed computing has been around for years
    - ✓ “distributed grid” is new, but grid computing is still computing science research
- 
-

# ***Distributed what exactly ?***

## ❖ Data

### ✓ Data Grids / Data Webs

- Remote data analysis and distributed data mining
- We don't need it today

## ❖ Compute

### ✓ CPU Grids

- Use compute power all over the world
- Programs must be “aware” of how to use distributed CPUs.

- **How many of those do we already have ?**

- **We should think GRID in the future developments, if necessary**



# *We have a cluster-GRID– we can distribute*

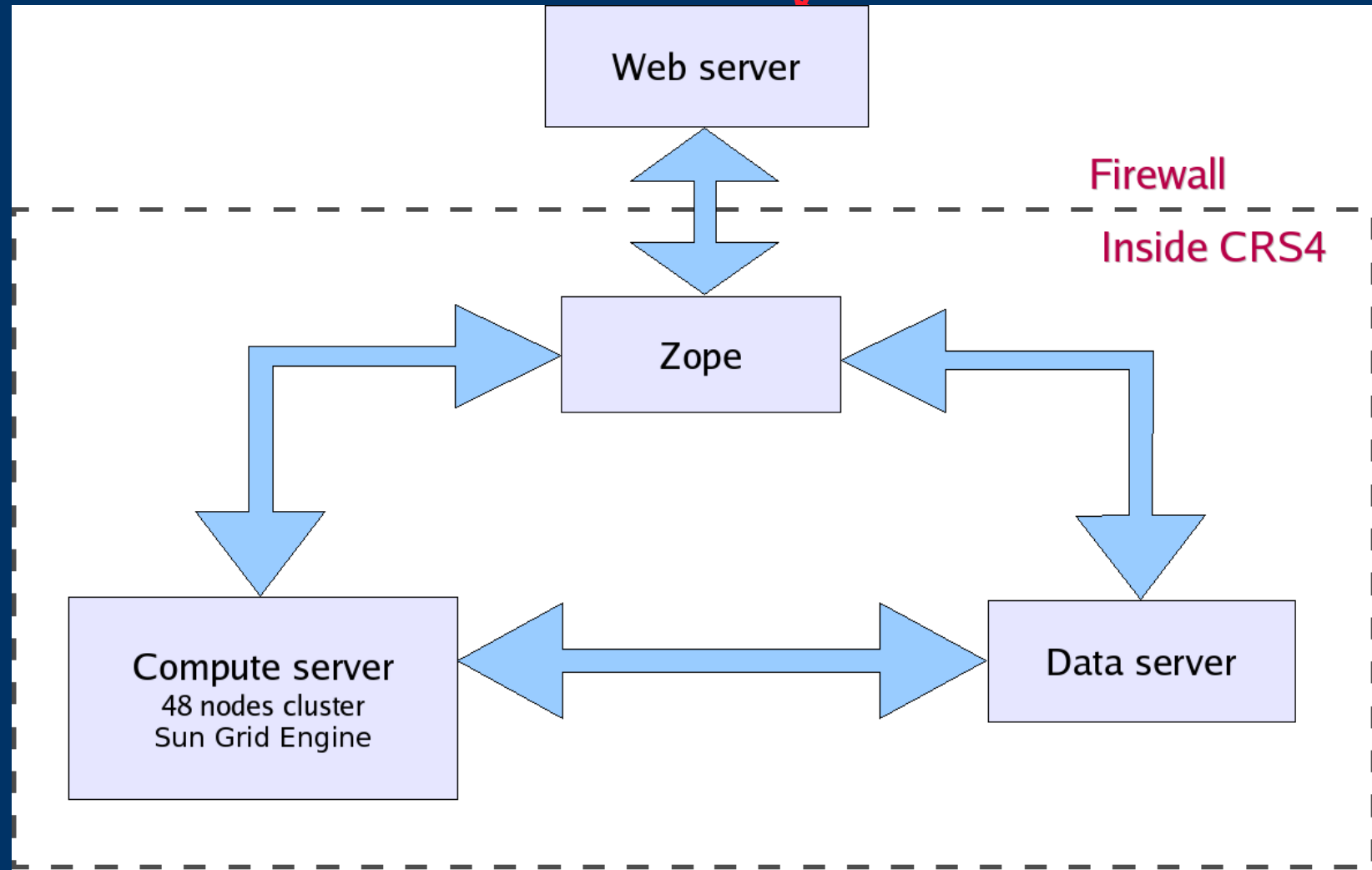
- ❖ Breaking problems into smaller pieces
    - ✓ Distribute pieces of data onto all nodes, distribute jobs on those nodes and merge the results when all jobs have completed
      - looks good for comparative genomics
  - ❖ Applications need to be parallelized to benefit from the grid
    - ✓ To “difficult” we leave that to the experts
      - (we are not interested)
  - ❖ Write new algorithms to solve new / existing problems currently unsolved
- 
-

# *First Problems encountered*

- ❖ Old programs were written on 32 bits machines
    - ✓ New processors: 64bits, Intel and AMD
    - ✓ Compilation and execution problems
      - “segmentation faults”
  - ❖ Make sure the results are identical and correct
    - ✓ This is true also when parallelizing a code or breaking it into pieces.
    - ✓ Because one should never blindly trust a computer...
  - ❖ Cost of migration (development)
- 
-

# Architecture

<http://www.bioinformatics.crs4.org>



# *Data server current setup*

- ❖ Target files like databases for blast or fasta are made available from the data server
  - ❖ Currently the directories are accessed as dynamic mount points -eg all nodes from the compute server can mount the data server bioinformatica disks
  - ❖ Why ?
    - ✓ our current jobs are short to medium size
    - ✓ time cost of copying each time the files is to high vs execution time
    - ✓ we run these jobs MANY times
- 
-

# *Future setup*

- ❖ Nodes specialized in short and medium size computation
  - ✓ dynamic mount of the data server directories
- ❖ Nodes specialised in long computations
  - ✓ copy of the necessary target files
- ❖ Nodes are already specialized,
  - ✓ but they currently all mount the directories





# *Surprise !!!!*

- ❖ yes, the GRID is raw computing power
    - ✓ but there is a ticket to pay to enter
  - ❖ it is slower when we send the calculations to the cluster
    - ✓ Delay between the job submission to the cluster and the time it starts executing
    - ✓ This delay may become longer than the waiting time on standard batch queues system
      - or even longer than the actual execution time of the job!
    - ✓ The cluster is best suited for long jobs, not for short ones.
- 
-

# Rules

- ❖ A set of rules are integrated into the application to either keep the job on the local machine or send it to the best suited SGE queue..
- ❖ ToDo: set of rules for checking the CPU usage of the Zope server and decide where to send the job
  - ✓ waiting for the production Zope machine

## *And still to do*

- ❖ the genotyping lab will produce sensitive data that should not be seen by other users
    - ✓ from the compagnies
    - ✓ related to patients
  - ❖ The CRS4 IT group is preparing for us a private subnet, which will be secured
    - ✓ data distribution
    - ✓ job execution
    - ✓ subnet will be dynamically generated when the job is submitted from a specific IP addresses range.
- 
-

# Workflows

- ❖ used to define processes for large-scale analysis
    - ✓ specifies what analysis need to be executed
    - ✓ the data flow between them
    - ✓ and relevant execution details
  - ❖ graphical workflow managers like Taverna help
  - ❖ but still, workflows are not easy to build
  - ❖ difficult for me, a seasoned bioinformatician... it will certainly be to difficult for our bioscientist users!
- 
-

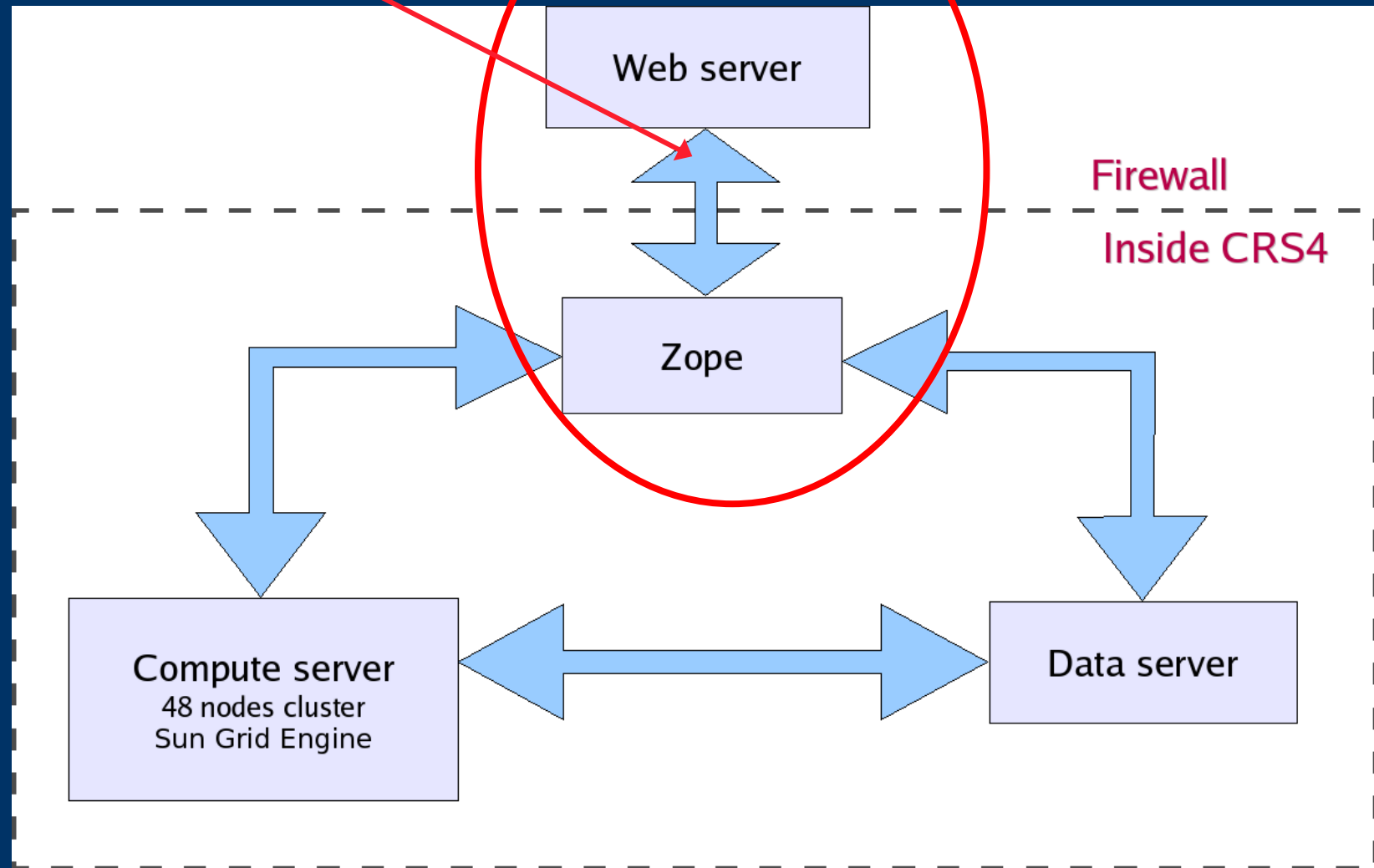
# *BioPortale*

- ❖ a central access point for all users
  - ❖ conceals the complexity of interacting with the Grid
  - ❖ provides a user friendly interface using Web form, which is a familiar sight for the user
  - ❖ defines static workflows
  - ❖ ToDo: more advanced facilities
    - ✓ like letting the user define its own workflow
    - ✓ but we need to educate them first!
      - bioinformatics-training the bioscientists is part of our mandate
- 
-

# Architecture

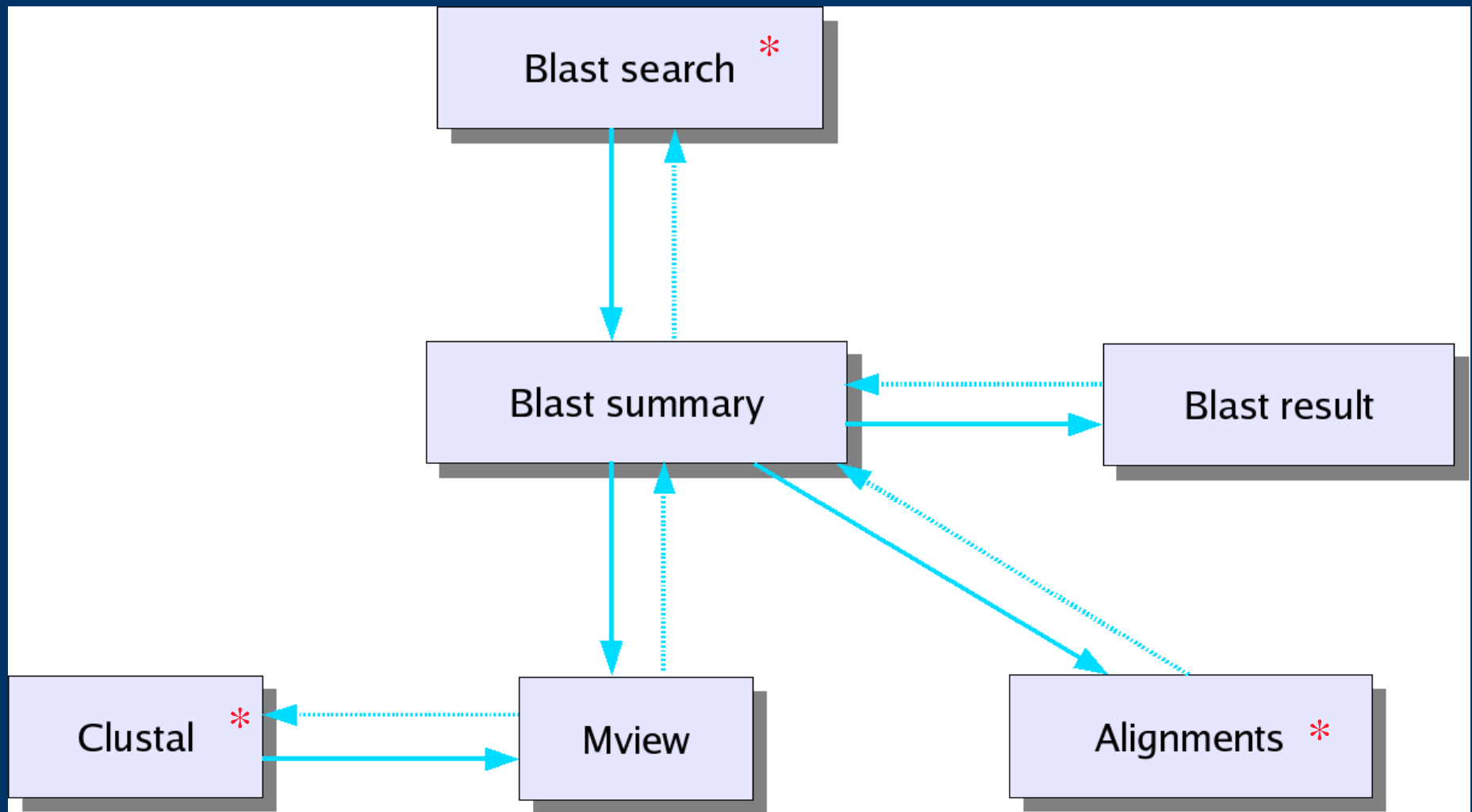
BioPortale

Proxy



# Our First Workflow

\* run on demand




Standard  
Blast form  
built as a  
Zope  
product  
and  
displayed  
by Plone

Blast — Bioinformatics Group - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

http://www.bioinformatica.crs4.org/analysis/Blast/blast

Bioinfo Main Page - BioInfoWiki Laposte.net web development MacBidouille.com : Bi... Perl modules documen... Submission - Papers

 Bioinformatics Group

Paste here your job ID  
  
Retrieve

home about us news events analysis modelling genotyping expression databases

Blast Search

DB type	Search title	DB name	Program
Protein	All the databases	UniRef100	blastp
Align views	Matrix	Exp.	Filter
pairwise	blosum62	default	false
Dropoff			default
Opengap	Extendedgap	Gapalign	Scores
default	default	true	default
Alignments			default

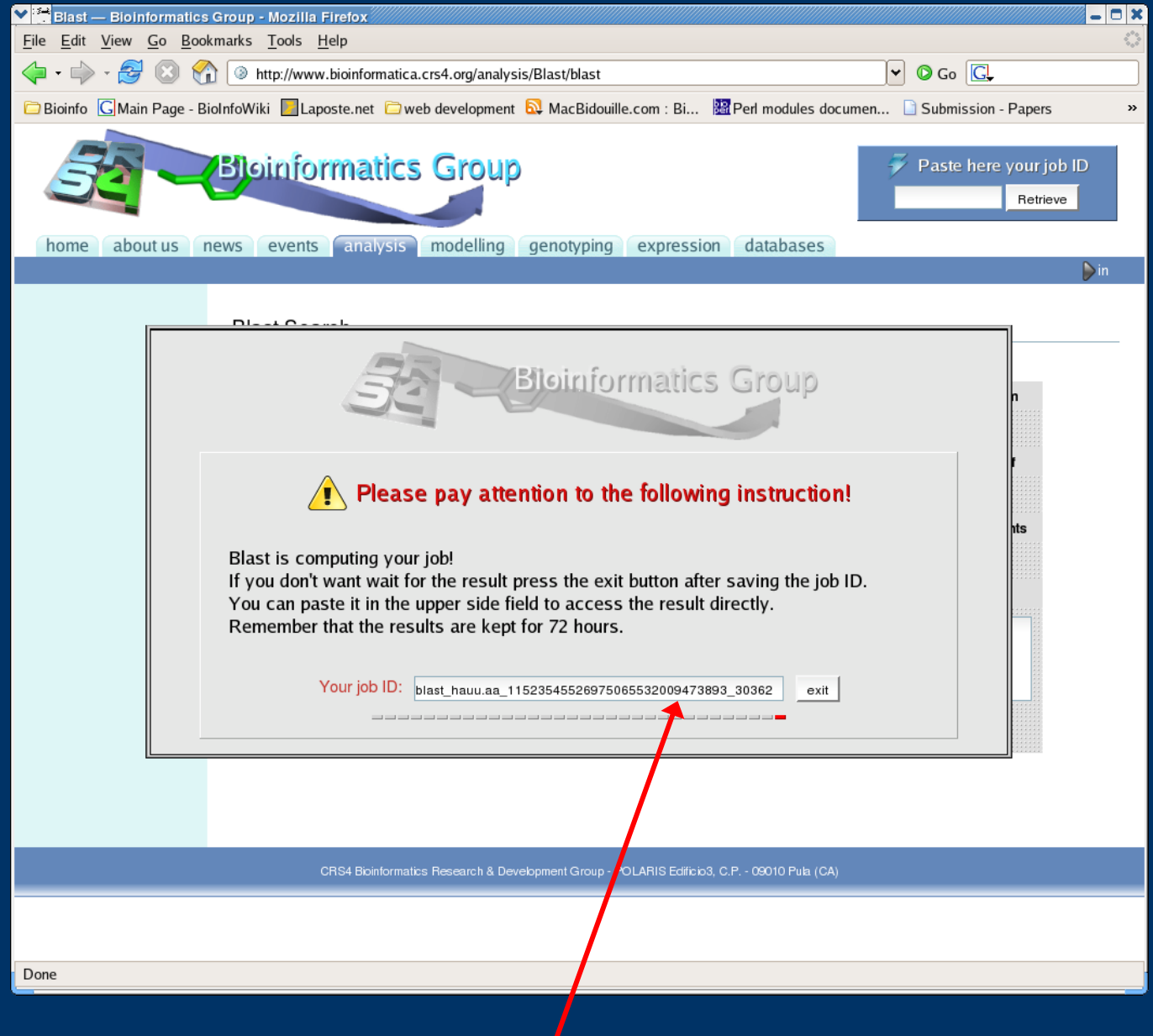
Paste a PROTEIN sequence or upload a file in fasta format.

CRS4 Bioinformatics Research & Development Group - POLARIS Edificio3, C.P. - 09010 Pula (CA)

http://www.bioinformatica.crs4.org/



The JobID is a 'session ID' that uniquely identifies a job. It is randomly generated, thus difficult to guess. The user can retrieve its job output at a later time;  
ToDo: retrieve the job status for long running jobs.



Job Id

Enter the job Id to retrieve previous results

ToDo: list the currently available jobsId for this user and their status

The screenshot shows the Blast search interface on the Bioinformatics Group website. The browser window title is "Blast — Bioinformatics Group - Mozilla Firefox". The address bar shows the URL "http://www.bioinformatica.crs4.org/analysis/Blast/blast". The page features a navigation bar with links: home, about us, news, events, analysis, modelling, genotyping, expression, and databases. A sidebar on the left is highlighted in light blue. The main search area is titled "Blast Search" and contains a form with the following fields:

DB type	Search title	DB name	Program
Protein	All the databases	UniRef100	blastp

Align views	Matrix	Exp.	Filter	Dropoff
pairwise	blosun62	default	false	default

Opengap	Extendedgap	Gapalign	Scores	Alignments
default	default	true	default	default

Paste a PROTEIN sequence or upload a file in fasta format.

Below the text input field, there are three buttons: "Browse...", "Run Blast", and "Reset".

At the bottom of the page, the footer text reads: "CRS4 Bioinformatics Research & Development Group - POLARIS Edificio3, C.P. - 09010 Pula (CA)".

The browser status bar at the bottom shows the URL "http://www.bioinformatica.crs4.org/".

points for  
static  
workflow  
access

Blast — Bioinformatics Group - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

http://www.bioinformatica.crs4.org/BRDG/analysis/Blast/result

Bioinfo Main Page - BioInfoWiki Laposte.net web development MacBidouille.com : Bi... Perl modules documen... Submission - Papers

**Bioinformatics Group**

home about us news events analysis modelling genotyping expression databases

Paste here your job ID  
Retrieve

### Blast summary

BLAST RESULT MVIEW SHOW ALIGNMENTS VIEW GRAPH CLUSTAL SUBMIT ANOTHER JOB ACTION TOOLS

Database	sgt	Nbr Entries	109903
Program	BLASTP	Version	2.2.13 [Nov-27-2005]
Gap extension penalty	2	Open gap penalty	11
Matrix	BLOSUM62		

CHECK ALL CHECK e.g. "1-2,5,7-9" SELECTION TOOL

1	Query	CAB96905.1 (gi 8980337 emb CAB96905.1)	Query length	1300
	Query description	cystic fibrosis transmembrane conductance regulator [Takifugu rubripes]		
	Number of hits	15		
	Alignment	DB:Name	Description	Length Score Evalue
1	☐	SGT:Y75B8A.26	SECSG##### Cloned	1144 168 6e-41
2	☐	SGT:Lmaj00188AAA	##### Cloned	1824 124 2e-27
3	☐	SGT:W09D6.6	SECSG##### Cloned	801 92 7e-18
4	☐	SGT:282163	JCSG##### Expressed	577 84 2e-15
5	☐	SGT:Rv1272c	U.Californ# Cloned	631 72 6e-12
6	☐	SGT:Tcru008725AAA	##### Expressed	388 69 5e-11
7	☐	SGT:282164	JCSG##### Expressed	598 68 9e-11
8	☐	SGT:Rv1273c	##### Cloned	582 67 2e-10

Done

# *More workflows*

- ❖ just replace Blast by Fasta, ClustalW, Muscle ...
  - ❖ the Zope products are modular and can be reused
    - ✓ it is just very long to develop products that can be reused
    - ✓ and difficult to find the people who have the know-how
    - ✓ (it is much easier to develop ad-hoc templates ... future legacy stuff)
  - ❖ these are the basis for the future workflows
- 
-

# *Some general thoughts after 4 months of development*

- ❖ Large number of “more suited” machines is better than a few “big ones”
    - ✓ Easier to upgrade
    - ✓ Redundancy and failover
  - ❖ Grid is better suited for long jobs than short ones
  - ❖ Grid is Not easy – not for the bench bioscientist
    - ✓ Need of friendly interfaces
  - ❖ Grid is still young, will evolve
    - ✓ get more or less easy ?
  - ❖ BUT
    - ✓ there is still the need for a powerful machine as web server to run the short calculations
    - ✓ and for large shared memory machines for some calculations or databases
- 
-

# *For what should we use the grid ?*

- ❖ to run many processes simultaneously
    - ✓ high number of jobs, example from the BioPortale
      - send jobs to different nodes of the grid
    - ✓ users that have many problems
      - send each problem to a different node
    - ✓ workflows that can be processed simultaneously
  - ❖ to run multi-process jobs for our internal research and collaborations
    - ✓ docking most certainly
    - ✓ genome comparisons
    - ✓ microarray analysis
    - ✓ image processing
    - ✓ genotyping
- 
-

# *The People*

- ❖ Paolo Zanella, Anna Tramontano, Patricia Rodriguez-Tomé
  - ❖ Giuliana Brunetti, Simone Carcangiu, Matteo Floris, Lisa Marras, Joël Masciocchi, Betta Muscas, Massimiliano Orsini, Enrico Pieroni, Frédéric Reinier, Alphonse Thanaraj, Maria Valentini  
(and new people soon)
  - ❖ And help from CRS4 IT group:
    - ✓ Lidia Leoni, Antonio Concas, Marco Pinna, Matteo Vocale, Carlo Podda, Alan Scheinine
- 
-



*Please visit us*

Parco Polaris  
Edificio 3

