



Databases and Functional Genomics Applications

CNR-ITB, INFN, CNRS, CILEA

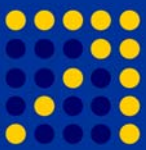


dkfz.

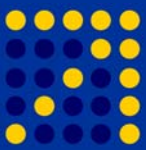


 cilea





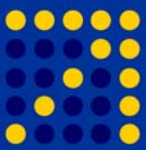
- **Biological Databases in Grid**
- Objectives
 - Provide availability of main Biological Databases in the Grid Environment
 - Optimize user wait times, and storage and transfer costs
 - Keep databases updated
 - Handle versioning
 - co-existence of multiple versions
 - handle naming, name clashes
 - store versions in a cost-efficient manner



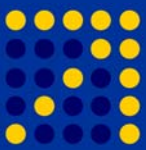
- **Functional Analogous Finder**
- **Challenge:**
 - find functional analogous gene products
- **No sequence comparison for this challenge:**
 - Small sequence changes sometimes produce big changes in the gene function
 - Genes compared according to their functional description in GO (Gene Ontology)
 - Statistical method: searched for possible functional analogous gene products by comparing the corresponding associated GO terms and their semantic similarity measure
- **Too expensive to be performed locally**
 - CPU time: 55 CPU-years



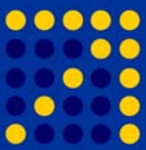
- **Testing of Biological Databases and Applications in GRID**
 - Testing Availability
 - Resilience: ensure usability and availability independent from failure of a single grid node or the central maintenance node
 - Ensure presence of the databases needed by the biologists
 - Testing Performance and Scalability
 - Applications scalability vs databases replicas number
 - More precise evaluation of database replication costs vs applications performance benefits. Tuning parameters for the update and replication engines.
 - Evaluate possible bandwidth optimizations
 - Evaluate / realize improvements



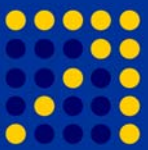
- Issue: Maintenance
 - Solution: timer based scripts polling for new versions of maintained databases from FTP sites
- Issue: Scalability / replication
 - 100-1000MB - sized files.
 - High replication costs.
 - However, not possible to use CEs not near to replicas of the needed database (remote download too time consuming)
 - unknown a-priori usage of each database
 - Useless to replicate rarely-used databases



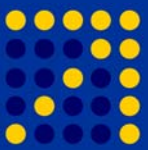
- **Best Evaluated Solution: Adaptive replication engine**
 - Database usage monitoring:
 - Moving average over N days (usually 10) for CPU-hours spent on every database
 - Configurable parameter for daily storage cost
 - Configurable parameter for cost of one minute of wait time for the user
 - Dynamic optimal replica number computation
 - math-based simulation of approximate user wait time with (N+1) and (N-1) replicas, based on the current wait time (N replicas)
 - Total (storage + user wait cost) cost computation with (N), (N+1) and (N-1) replicas. Iterative algorithm to find optimal number of r.
 - Constant adaptation of the number of replicas



- Issue: versioning
 - Old versions of databases can be meaningful to biologists
 - Need to replicate computations /demonstrate previous results obtained
 - However: high storage costs.
 - Versions updated every few months.
 - In addition: risk of name clashes
 - Need to be able to specify one precise version, unambiguously
 - Co-existence of files with same name



- Solution: (versioning)
 - Storage costs: patches
 - Only xdelta patches are stored, and only once (no replication)
 - One patch regresses one version
 - Patches are small, they ideally represent only the differences between two versions of a file.
 - Patches can be applied in sequence to regress to any earlier version
 - Every time a new DB version $V+1$ is released, the DB version $V+1$ is uploaded, DB version V is deleted, and a patch $V+1$ to V is uploaded.



- Solution: (versioning)
 - Name clashes:
 - Names for older versions are encoded with a specific convention
 - Python script “preparedb.py” needs to be used on the Worker Node for downloading an older version of a required database. The version is given as a date in ISO standard YYYYMMDD. Download of patches and regression to the required version is performed automatically



- Preliminary testing and considerations

Size of databases

- Max=5.3GB (nt.gz); **Avg=691MB**; $\sigma / \bar{x} = 1.59$
- Replication speed (bandwidth testing)
 - Avg=1.91MB/sec; $\sigma = 1.02$; $\sigma / \bar{x} = 0.53$
- Replication balance (hours/day n-replicas equilibrium)
 - Up to now we have tuned the parameters so that:
 - 1GB file will get 1 additional replica for each 10 CPU-hours/day of average computation (moving average over 10 days).
 - 2GB files get half that number of replicas, etc.
 - We will perform more exact cost computations/optimization



Biological Databases - Use case

BioinfoGRID

- Biologist needs to blast 50,000 sequences against NR version dated 2005-06-15
 - Blastp CPU time against NR: 25 seconds per sequence
Total: 15 CPU days (Xeon class CPU)
- How is that with the Grid
 - Blast is a widely used application so we support in various ways:
 - Through the BGBlast (BioinfoGridBlast) application
 - Through a more low-level approach which can be used for any kind of application (not just Blast)
 - Through the Gilda portal



- Using BGBlast (BioinfoGridBlast)
 - From the commandline on an UI w/ bgblast installed
 - `bgblast.pl -p blastp -i input.seq -d nr --version 2005-06-15`
 - this will take care of
 - splitting the input.seq reasonably
 - launching parallel jobs
 - regressing the database on the WN to 2005-06-15
 - waiting for the results
 - Needs the commandline: a few users might find uncomfortable



- Using the low-level approach (this works with any application, not just Blast or widely known ones)
 - Inside the job, launch: `./preparedb.py nr.tgz 2005-06-15`
 - This will download the requested (NR) database
 - Additionally, it will regress it to the (optionally) specified date by downloading and applying the needed patches
 - Submit the job to the grid. Bind execution near a replica by using: `edg-job-submit --resource someCE .` Place `preparedb.py` in your sandbox
 - Monitoring job completion and fetching results is your responsibility



- Using the Gilda portal (Multi-Blast)
 - This takes care of:
 - Launching Blast (no database regression at this time)
 - Graphically monitoring your jobs' status
 - Intuitively showing all the steps in a graphical user interface (web)

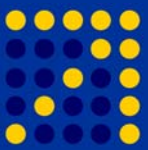
Biological Databases - Use case

BioinfoGRID

The screenshot shows a Mozilla Firefox browser window displaying the GENIUS Grid Portal. The address bar shows the URL <https://glite-demo.ct.infn.it/>. The page features a navigation menu with links like 'Getting Started', 'Latest Headlines', 'GILDA', and 'EUCHENA'. Logos for INFN, eGEE, genius, and NICE are visible. The main content area is titled 'Submit MULTI BLAST' and contains the following form fields:

- Database Name:
- Database Version:
- Database Type:
- Input query:

Below the form is a button labeled 'Configure the dataset to submit your MULTI Blast to the GRID'. The left sidebar shows a tree view of services under 'BLAST Services!', with 'Submit MULTI Blast' selected. The footer contains copyright information: 'Copyright © 1998 - 2006 Nice S.r.l. All trademarks and logos on this page are owned by NICE s.r.l. or by their respective owners.'



Biological Databases Available in GRID

- The following biological databases are currently available in Grid:
 - InterPro databases:
 - SUPERFAMILY [98MB]
 - TIGRFAMs (Haft, D.H. et al. 2001) [36MB]
 - Pfam (Bateman, A. et al. 2000) [90MB]
 - PROSITE Prosite.Patterns (Hofmann, K. et al. 1999) [36KB]
 - PROSITE Confirm.Patterns (Hofmann, K. et al. 1999) [49KB]
 - PROSITE profile (Hofmann, K. et al. 1999) [1.5MB]
 - PRINTS (Attwood, T. K. et al. 2000) [7.4MB]
 - PRODOM (Corpet, F. et al. 1999) [5.6MB]
 - SMART (Schultz, J. et al. 2000) [5.4MB]
 - PIRSF [1.7MB]
 - PANTHER [1.1MB]



– BLAST databases:

- nr (NCBI) [0.98GB]
- nt (NCBI) [5.0GB]
- pdbaa (NCBI) [5.8MB]
- UCSC_human_chrs (UCSC) [899MB]
- human_genomic (NCBI) [1.62GB]
- refseq_protein (NCBI) [720MB]
- refseq_rna (NCBI) [360MB]
- refseq_genomic (NCBI) [2.02GB]
- ecoli (NCBI) [906KB]
- yeast (NCBI) [1.86MB]
- uniprot_sprot (UNIPROT) [41MB]
- uniprot_trembl (UNIPROT) [631MB]
- est_human (NCBI) [1.30MB]
- est_mouse (NCBI) [741MB]



Goal

Compare gene products according to their description
AND NOT
according to their sequence similarity.

As description we use the standardised terminology of the Gene Ontology (GO).

Data source:

The Gene Ontology Database (GODB), is a repository of the GO and the associations between the terms and the gene products (GOA).

Currently there are 2M gene products described by 21000 terms producing 9M associations.



Approach

- A selection of about 1M “well annotated” gene products are involved in the search.
- A simple chi-square application compares the common and non-common terms between two compared gene products.

Problem

- A comparison of one gene product against the whole 1M gene products occupies 1 CPU for 30 min. on average
- The whole “every gene product against each other” search would occupy 1 CPU for more than 50 years.



Solution

- Split the search into a number of small jobs and distribute them together with the DB on as many free WN as possible.

The **job submission** is made by a script **running as a daemon**

- The script submits 80 jobs every 30 minutes
- It is possible to **run more instances of the submission daemon** in order to **increase the total number** of jobs submitted in one hour
- The **multi-process** submission **improve the speed** of submission

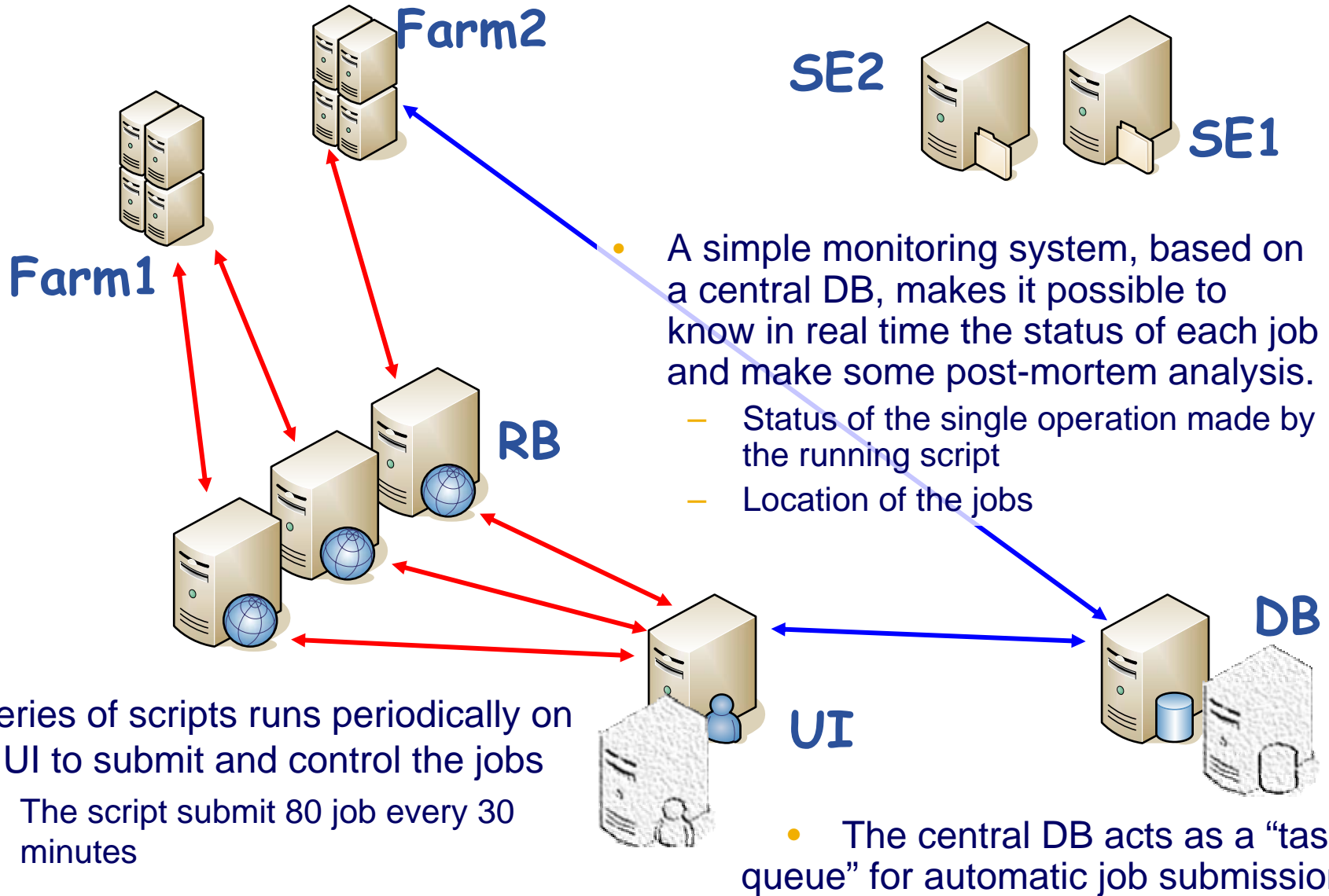
The submission uses **3 RB** in a round robin algorithm in order to **avoid overloading** a single RB and to avoid that the **failure** of a single RB can stop the submission of jobs

- Retrieve periodically the OutputSandbox of the jobs
- **Monitor** the status of the production by simply querying the monitoring DB
 - The user can know the **number of processed/running genes**
 - The **number** of the **running** jobs
 - The **location** of each job
 - **Debug** possible **errors** in running jobs
- The software to submit jobs is installed on 2 different machines in order to avoid that a single hardware failure can stop the submission

Functional Analogous Finder

the job submission

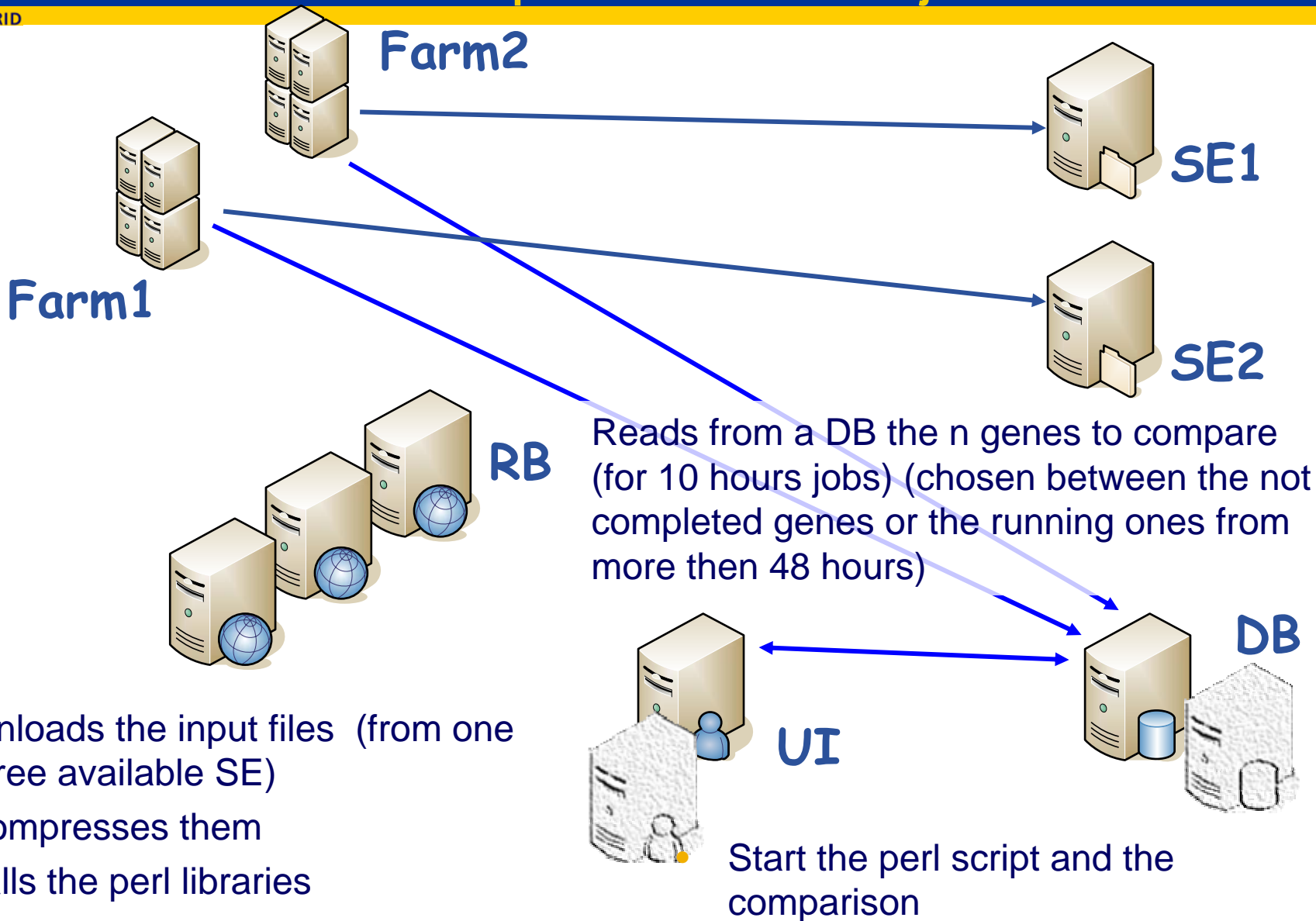
BioinfoGRID





Functional Analogous Finder

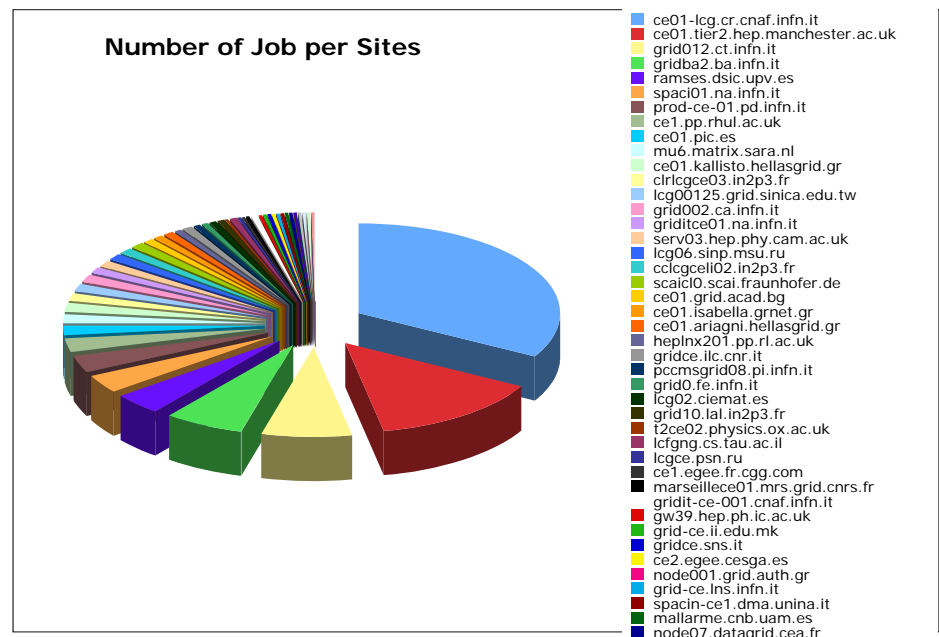
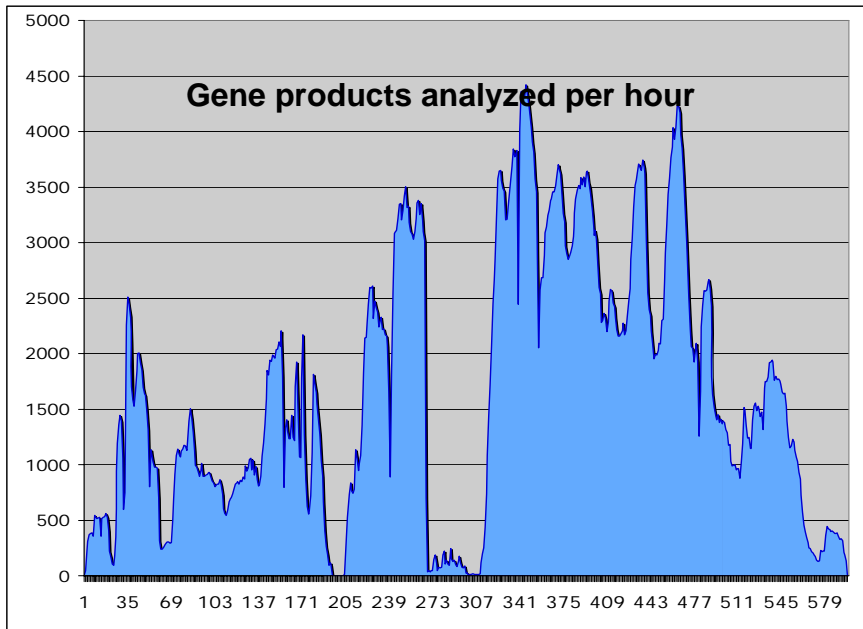
actions performed when the job reaches the WN



Functional Analogous Finder Results

BioinfoGRID

- All 1M gene products processed in **less than one month**
- **Different Farms used: 64**
 - **Different Hosts used: 2446**
 - Total Submitted jobs: 95041
 - Total started jobs: 66313
 - Total of successful jobs (from application point of view): 42992
 - Total Failed job (for input staging problem): 3209



- Good efficiency
- Easy to manage
- A quite general procedure for the submission of jobs
 - A general procedure to keep track of to-do/done tasks
- Grid submission success rate is not so important
- The procedure is thought in order to run for a long period in an unattended way
- The rate of analysis grows linearly with the number of nodes that are available
 - We don't see any bottleneck at this level (the task queue does scale on the thousand concurrent jobs)
- It is important to avoid black-holes (WN or farms that make job to fail always)
 - A static fix for black-holes has been made on the running script
 - There is the possibility to exclude farm from submission by modifying the jdl file
 - An automatic procedure is needed to create a dynamic black-list of farms

- We will set up a database (best hits DB) publicly accessible to query the search results and downloading the data.
- The GODB is frequently updated and therefore the functional analogous search needs to be re-run after every update.

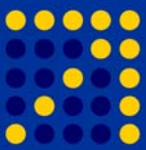
The full search lasted only 30 days, however, it is unthinkable to repeat the search at every updates of the GODB.

Therefore we are currently working on setting up a procedure that, in case of an update of the GODB, will limit the search to only those genes affected by the update, and makes the necessary adjustments in the “best hits DB” in order to reflect the update in the GODB.

- **The procedure will be improved using the features of the new gLite WMS**
 - **Pre and Post running script (shallow and deep re-submission)**
 - **File picking**
 - **Bulk (parametric) submission**



- Gabriele Trombetti, Alessandro Orro, Ivan Merelli, Luciano Milanese - “**BGBlast: A Blast Grid Implementation with Database self-Updating and Adaptive Replication**” - *NETTAB2006 proceedings*
- Milanese L, Andreas G, Arlandini C, Beltrame F, Bishop C, Breton V, Ernest P, Jacq N, Legre Y, Liò P, Liuni S, Mazzuccato M, Maggi G, Meloni G, Merelli I, Morra G, Orro A, Porro I, Sanger M, Shuai S, Trombetti G - “**BioinfoGRID: Bioinformatics Grid Application for life science**” - *BITS2006 proceedings*
- Alessandro Orro, Ivan Merelli, Gabriele Trombetti, Luciano Milanese - “**Enabling Post Processing Data Extraction in Grid Environment**” - *NETTAB2006 proceedings*
- Luciano Milanese, Claudio Arlandini, Paride Dagna, John Hatton, Mario Amdrea Marchisio, Andrea Mattasoglio, Ivan Merelli, Alessandro Orro, Marco Pirola, Paolo Ramieri, Giampietro Tecchiolli, Gabriele Trombetti, Pierfrancesco Zuccato - “**BioinfoGRID: Bioinformatics GRID based applications overview**” - *NETTAB2006 proceedings*



- Gabriele Trombetti
- Alessandro Orro
- Ivan Merelli
- Andreas Gisel
- Giacinto Donvito
- Giorgio Maggi
- Claudio Arlandini
- Luciano Milanese
- Giuseppe La Rocca
- Roberto Barbera

