



# A Grid Approach for Large Data Processing in Biomedicine

*Fato Marco, Papadimitropoulos Adam, Porro Ivan,  
Scaglione Silvia, Schenone Andrea, Torterolo Livia, Viti  
Federica*

Department of Communication, Computer  
and System Science (DIST)  
University of Genoa, Italy  
*pivan@dist.unige.it*

- A grid platform built with scalability and interoperability in mind
- It will allow large scale (“thousands”) bioinformatics data analysis and will have a simple web-based interface
- Implemented for microarray (first case study) but designed to be flexible enough to handle different data types



# *What not*

- It's NOT a Grid-based database
- It's NOT (yet!) a tool to access everything everywhere
- It's NOT a closed solution (in terms of licenses and applications)
- It's NOT designed to be a vertical application



# Why

- Microarray have intrinsic costs (chip, raw materials, manpower), about 500€ each chip
  - laboratories usually work on small experiment (10 – 30 uA)
  - the better statistical analysis you want, the larger the dataset must be ( $N^{\circ}$  samples <<  $N^{\circ}$  variables)
  - client side software tools interfere with data share
- Bio"everything" data is growing in size and speed of acquisition
  - Data mining is difficult
  - Data storage is challenging (hierarchical vs. on-line tradeoff)

- Identify free open source tools for microarray analysis  
→ **dChip**
- Identify a network infrastructure to provide secure and reliable storage → **gLite Grid**
- Design a modular system, open to third party modules (algorithms, workflows) → **W3C standards**
- Implement a user friendly web interface prototype
  - Available: **Genius**
  - Custom: **PHP + Apache, J2EE, portlet**



# *dChip porting*

- dChip is a Windows GUI application for microarray data analysis
- One of available options, choosed as a case study
- **Q:** Is it possible to run dChip command line on a Linux environment and, maybe, over The Grid?
- **A:** Yes. Steps:
  - Linux port
  - Investigate parameters and configuration
  - Standardize I/O files (also intermediate ones)
  - Separate the program into standalone modules
  - Make them parallel

**Options**

Clustering | Analysis | Model | Chromosome |

**Preprocessing and algorithm**

- Standardize rows (subtract Mean and divide by SD)
- Pre-calculate distances
- Distance metric: 1 - Correlation
- Linkage method: Centroid
- Gene ordering: By cluster tightness

**Visualization**

- Red/black/green coloring
- Sample names always visible
- Averaged gene profile pattern
- Add new color for Control+Click
- Show probe set name

Displaying range of standardized values: 3

Number of letters shown for sample information: 1

P value threshold for calling significant clusters

Gene:	0.001	Sample:	0.05
-------	-------	---------	------

**Options**

Clustering | Analysis | Model | Chromosome |

Curve along chromosome

- Min:
- Max:
- Threshold:
- Use Min and Max as threshold
- Linkage curve file:

LOH analysis

- Show only first sample name
- Hide conflict LOH
- Inferred LOH method: Hidden Markov Model
- Inferred LOH call threshold: 0.5
- Reference genes:
- Remove LOH regions consistent with 10 % of

Copy number analysis

- Use paired normal as reference
- % of samples trimmed:
- Parent-specific copy number
- Truncate small loci
- Inferred copy method: Median smoothing
- Inferred copy file:
- Inferred copy step: 1

SNP marker

- Use SNP-specific frequency
- HMM length: 100
- Genotyping error: 0.01
- Average Het rate:

**Options**

Clustering | Analysis | Model | Chromosome |

**Open group**

- Log transform expression values
- Show "R View" icon
- Allow TXT/info file to have unknown probe sets
- Working directory: C:\Documents and Settings\Luca\Desktop
- Search and save DCP files in the Working Directory
- Load probe data in memory
- Array only has PM probes

**Output**

- Additional gene information (e.g. Accession, Affymetrix description)
- Insert Excel and Image outputs into the Analysis View
- Treat array outlier as missing expression value (also in clustering)
- Gene name from LocusLink when available (also in clustering)
- Gene list file

  - Mask redundant probe sets
  - Omit Affymetrix control probe set at filtering and comparison

- Show online link dialog when accessing Internet
- Consider measurement error when averaging

**Model-based expression value**

Method: PM/MM difference model

- Check single outlier
- Check probe outlier
- Check array outlier
- Truncate negative PM/MM difference to
- Treat image spikes as single outlier
- Do not call all replicate arrays as array outlier

Exclude: 0 5' probes

**Probe sensitivity index (PSI) file**

Usage: Do not use

File: C:\Documents and Settings\Luca\Desktop\dChip

**SNP array**

Compute signals separately for A and B allele

**Normalization method:** Invariant set normalization

**samples (after pooling replicate arrays):**

standard deviation / Mean < 1000

arrays used >= 20 %

replicate arrays called Present

median(Standard deviation / Mean) < 0.5

level is >= 20 in >= 50 % samples

**all genes**

documents and Settings\Luca\Desktop\ make sure the file is closed

**ions...**

**Reset Default**

# *Computing & Data Grid: gLite*

- Secure and distributed data storage
- Access to several thousands CPU
- LCG (LHC Computing Grid) is the middleware born to support forthcoming (2007) High Energy Physics experiments
- gLite is its evolution (deployment in production on April 2006)

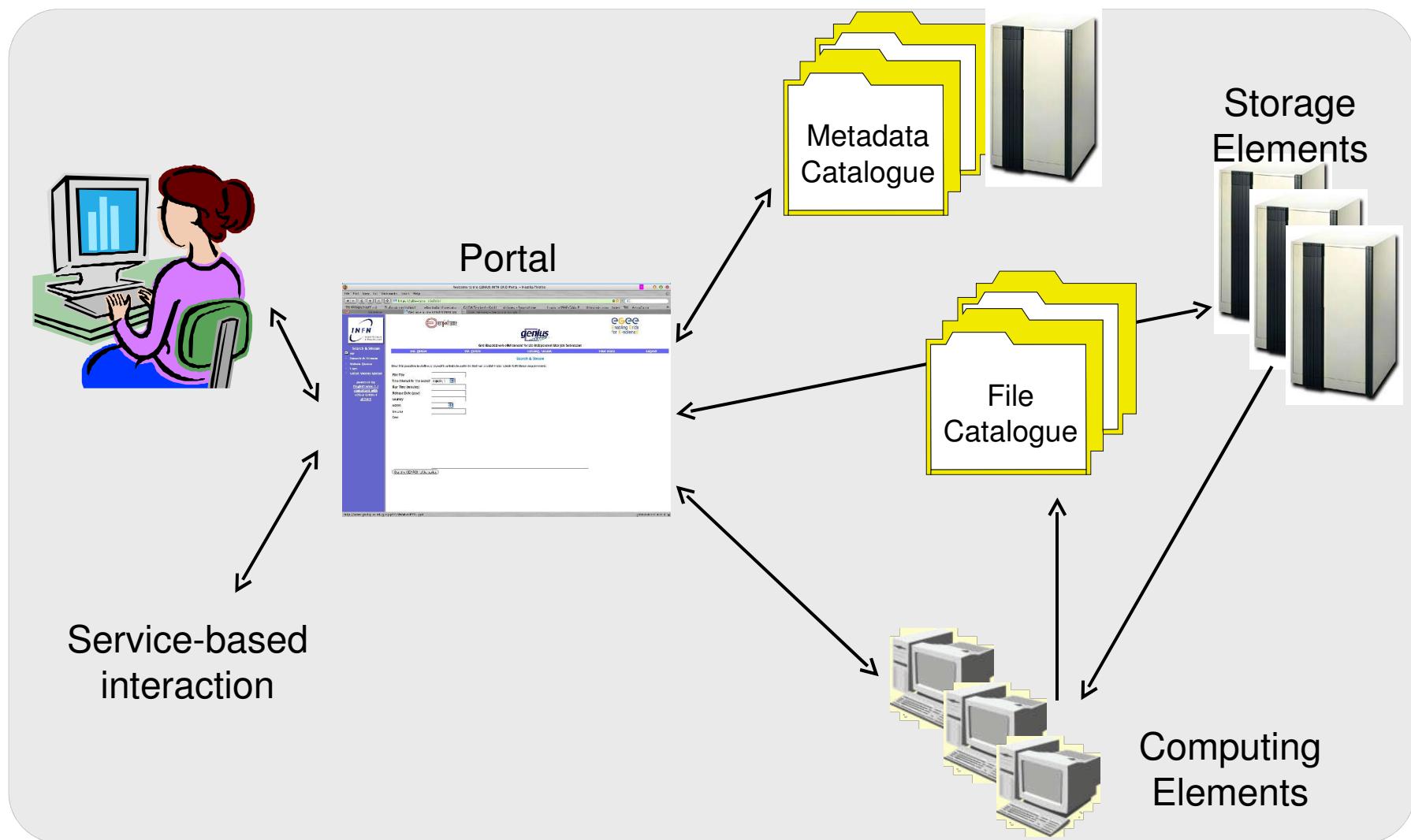
# Metadata: MIAME & AMGA

- Metadata are relationally structured information associated to experiment data
- MIAME (Minimum Information About Microarray Experiment) is a MGED de-facto standard
- We defined a custom ontology for stem cells experiments to maintain homogeneity in collected data among different laboratories
- AMGA (ARDA Metadata Grid Application) is the Metadata catalogue of gLite 1.5, with replication and federation support

# Interface prototype: Genius

- Provide a open source grid-portal based on Nice Srl product Enginframe
- Features:
  - Grid Security Infrastructure compliant
  - Modular and flexible: it can handle almost every programming and scripting language for underlying logic
- It could publish its interface components also as web services (WSDL)

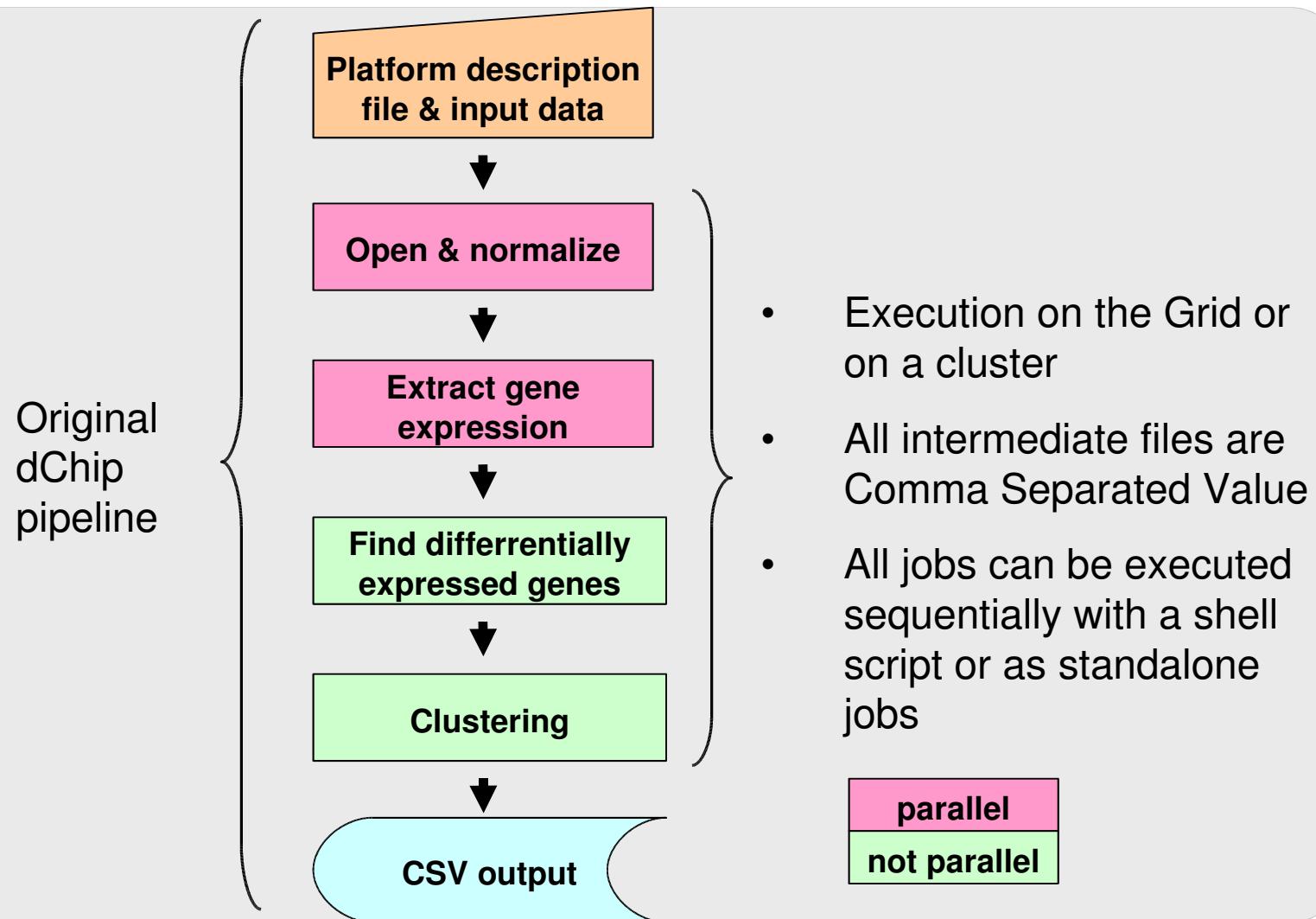
# System Architecture



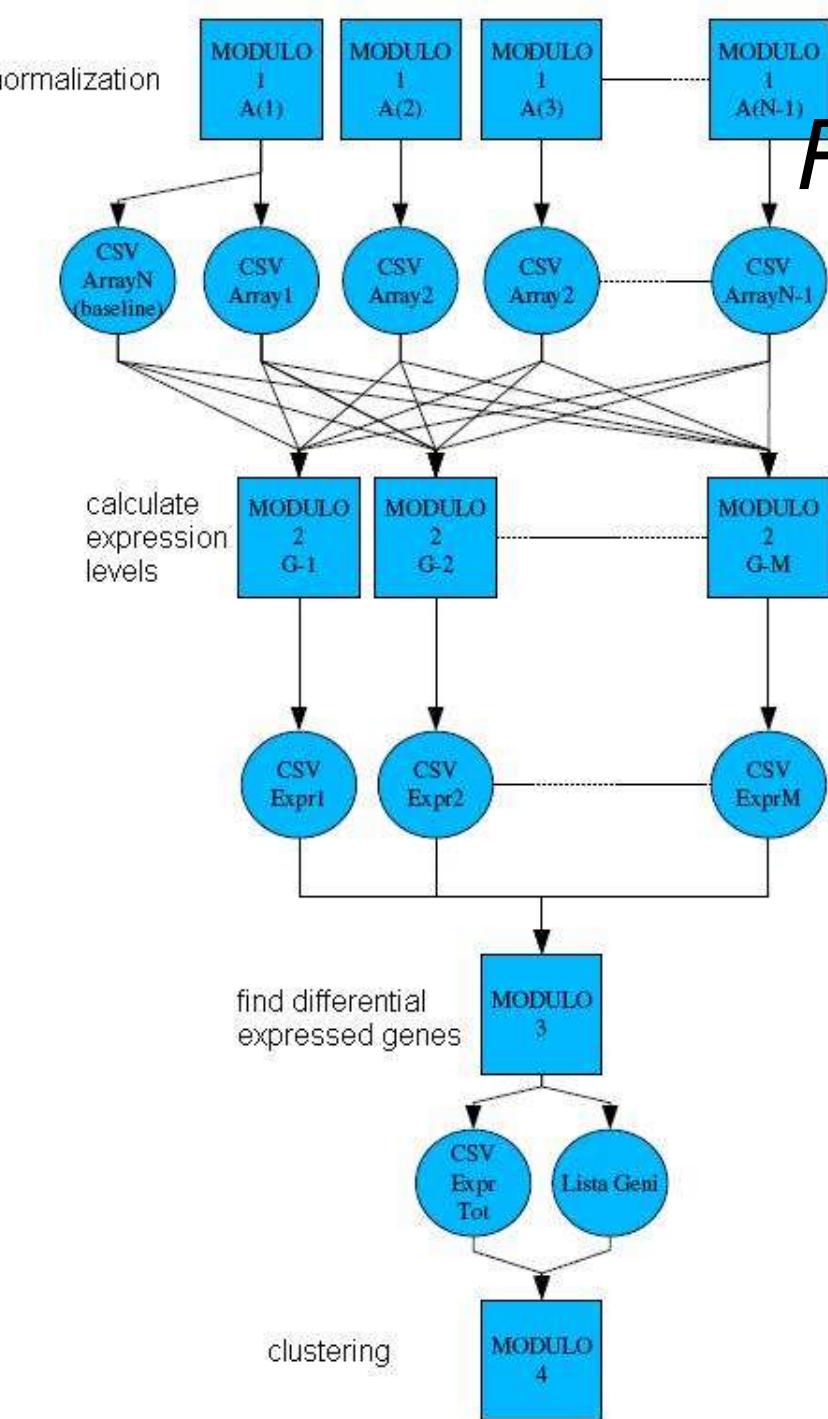
# *MPI-dChip available*

- dChip has been successfully ported to Grid and parallelized with MPI (mpich2)
- Two version available (Linux):
  - standard file I/O for MPI clusters (RHEL3/4)
  - grid GFAL I/O for MPI grid job (gLITE 3.0)
  - mix available
- Source code is GNU/GPL
- Availability: on request [LITBIO portal? Sourceforge?]

# MPI-dChip structure



# Parallelization schema



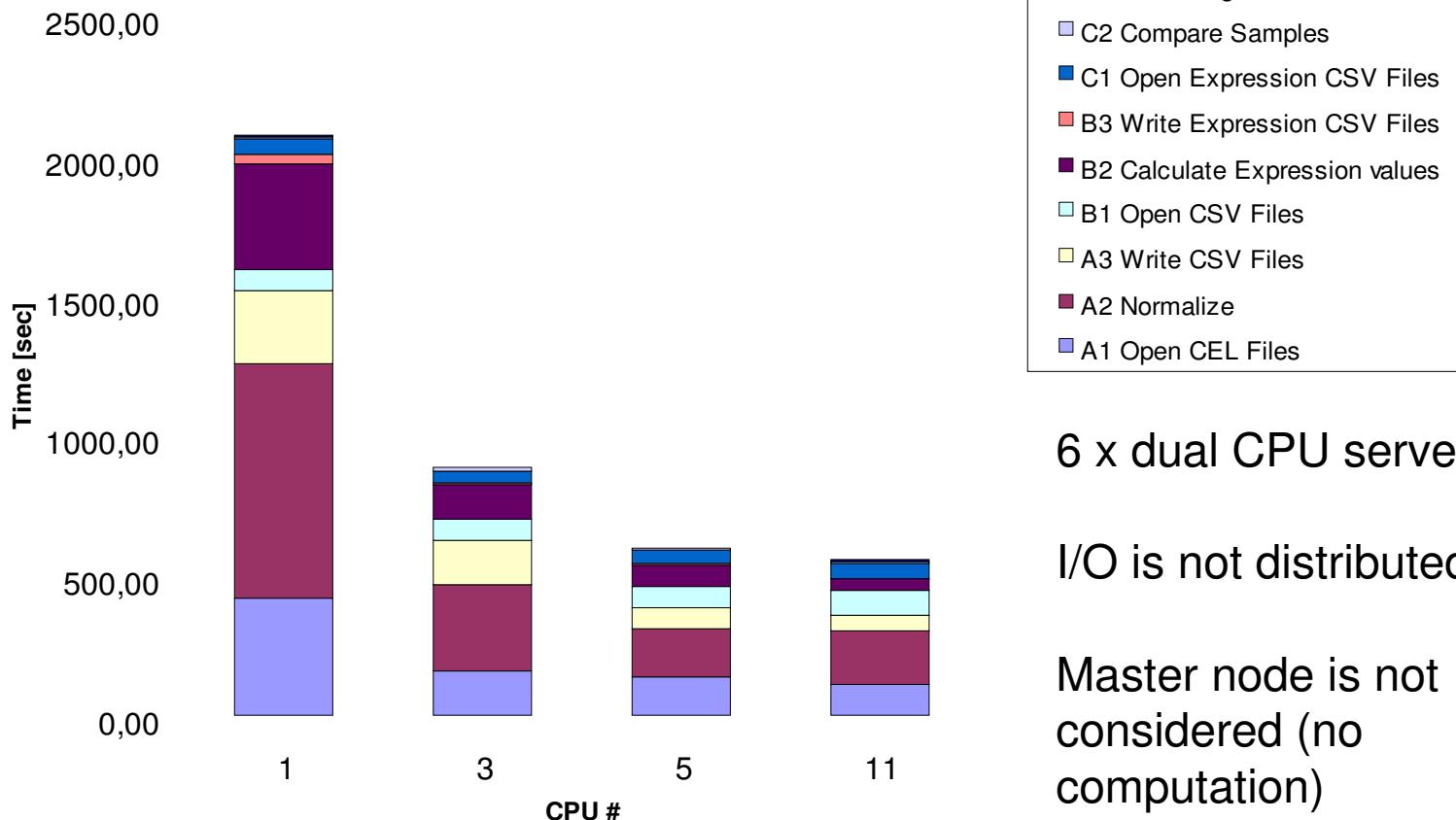
e.g.: we have 10 job and 200 files, each with 20.000 genes

Each normalize job takes 200/10 files and the base file and produce 200 normalized input files

Each model (PM or PM/MM) job takes 20.000/10 genes and produce 10 expression level files

Differential expressed genes calculation and clustering are not parallel. Results are saved to 2 CSV files

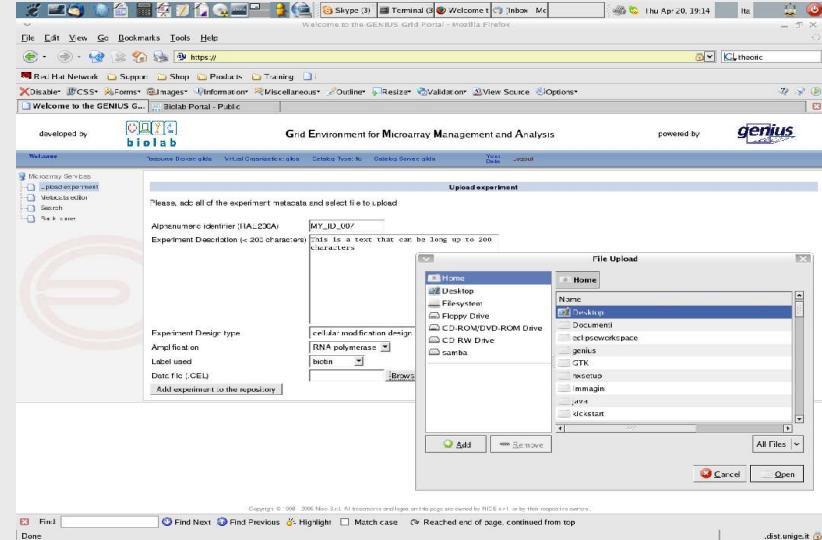
## dChip MPI tests



We successfully installed and customized Genius Grid portal

We developed microarray data management service prototype:

- Annotated experiment upload
- Search
- Edit



# *Open issues & future work*

- Genius 3.1 / Enginframe under extensive testing
  - Standard development is easy
  - Low level customization and software maintenance could be difficult
- Workflow enactment with Moteur engine partially successful
  - Moteur limitations for data and services
  - Genius integration has been undertaken
- Webservice available only for Genius interface services
- Integrate other tools
  - ArrayExpress Import/Export
  - R/Bioconductor based tools

# Acknowledgments

- This platform is currently part of the Italian FIRB project LITBIO (Laboratory for Interdisciplinary Technologies in BIOinformatics), <http://www.litbio.org>
- Thanks to Ulrich Pfeffer, IST Genoa for providing data for our tests and to Luca Corradi, thesis student, for coding MPI dChip
- Questions?