*École des sciences avancées de Luchon*
**School for advanced sciences of Luchon**

**Networks and data mining**
**Session II, June 27 - July 11, 2015**

# Integrating multi-omics

*Luciano Milanesi*

Consiglio Nazionale delle Ricerche

**InterOmics**
Flagship Project

- Introduction
- Omics challenges
- Data Integration
- Big Data
- Personalized system medicine
- International Initiatives
- Conclusions

The *"Omics Sciences"* consist of several areas of investigation :

- **Genomics,**
- **Proteomics,**
- **Interactomics,**
- **Bioinformatics,**
- **Neuroinformatics**
- **System Biology**
- **Metabolomics**
- Ecc.

These and the correlated disciplines constitute the paradigm around which all the research in the fields of biomedicine, biotechnology and ICT generally applicable to the biomedical sciences

# SNP and Biomarkers Analysis

InterOmics CNR

# Omics Data Explosion



EMBL-Bank growth

EMBL-Bank Growth
21-Nov-2011

# Interactomics and Pathways Discovery

Bioinformatics
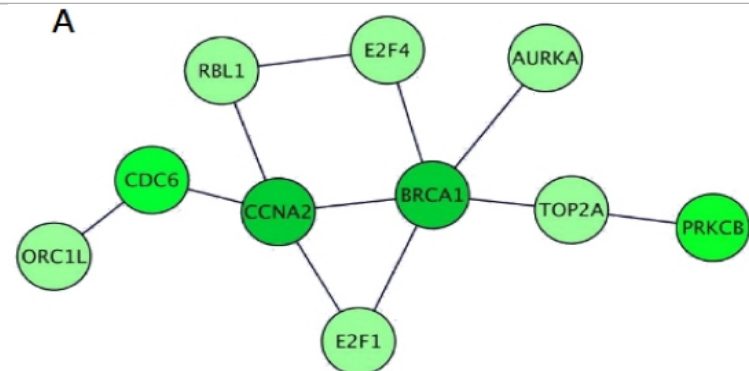
System Biology

Biotechnology

System Medicine

ICT

- Definition

    - **Big Data** refers to a collection of data sets so large and complex that it's impossible to process them with the usual databases and tools.

    - Because of its size and associated numbers, Big Data is hard to capture, store, search, share, analyze and visualize.

    - The three V's: **Volume, Velocity, Variety**

    - **High-Volume**: Amount of data

    - **High-Velocity**: Speed rate in collecting or acquiring or generating or processing of data

    - **High-Variety**: Different data type such as audio, video, image data, sequence data

- Processing

    - Parallel processing (eg. *Hadoop)*
    - Processing of data sets too large for transactional databases
    - Analyzing *interactions*, rather than *transactions*

# Big Data

Collection – get the data

Storage – keep the data

Querying – make sense of the data

Visualization – see the scientific value

**Medical Science**

- Data bases from
- e-Health
- Patient Records
- Medical ImagingMRI & CT scans,, …
- Telemedicine
- Genomics
- Environmental data
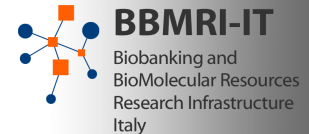- Food science
- Biosensors

**Big Pharmaceutical Companies**

- **Cloud computing** in combination with **Big Data Tools** can be used to obtain the power and the scale of computation required to facilitate large-scale efforts required in **translational medicine data integration** and to perform analysis in more efficient and economical way.

- ## Resources:

  - HPC (High Performance Computing) Cluster
  - HPSI (High Performace Storage Infrastructure) DDN –
  - WRVM (Web Remote Virtual Machine)
  - Databases: MySQL, ORACLE, SQL Server
  - Cluster Intel Servers: 44
  - Total RAM: 2.080 GB
  - Total Disk space: 1.164 TB
  - 192 CPU and 1.216 core
  - GPU Server : 16 GPU, 16 CPU and 96 core
  - Operating system: Ubuntu 13.04, Centos 6.5, Window Server, Mac OS
  - Portal technology: Java portal (LIFERAY)
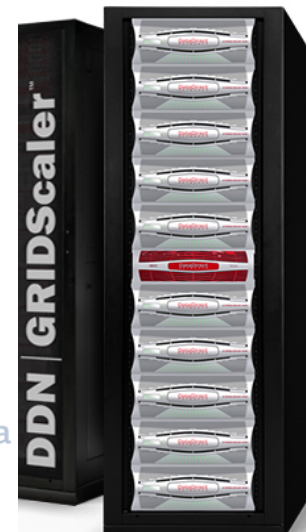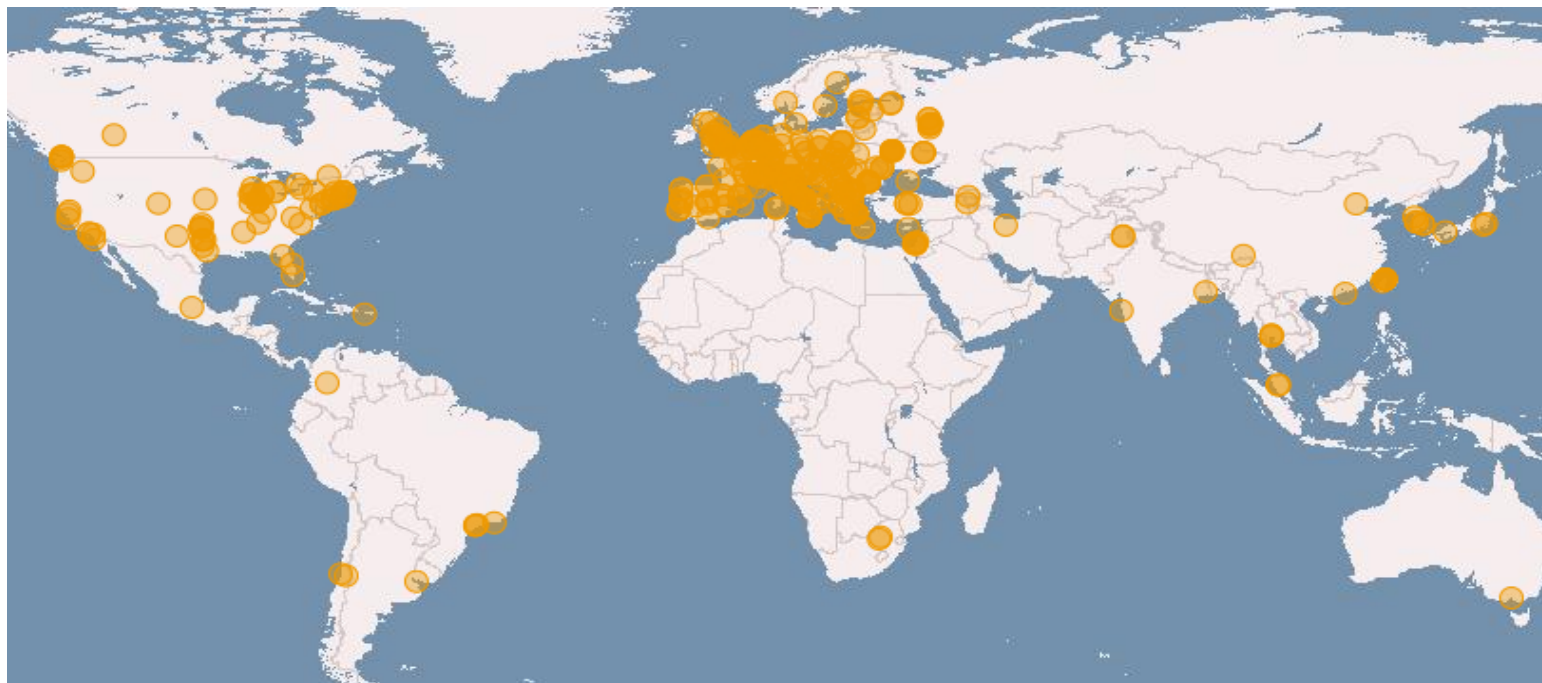  - GRID Node
  - Virtual Node
  - Cloud Computing
  - Hadoop



**BBMRI-IT**
Biobanking and
BioMolecular Resources
Research Infrastructure
Italy

**InterOmics**
Flagship Project

Ministero dell'Istruzione
dell'Università e Ricerca

$\hat{\sigma}^2$
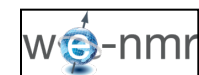**MIMOmics**

**HIRMA** Hepatocarcinoma
Innovative
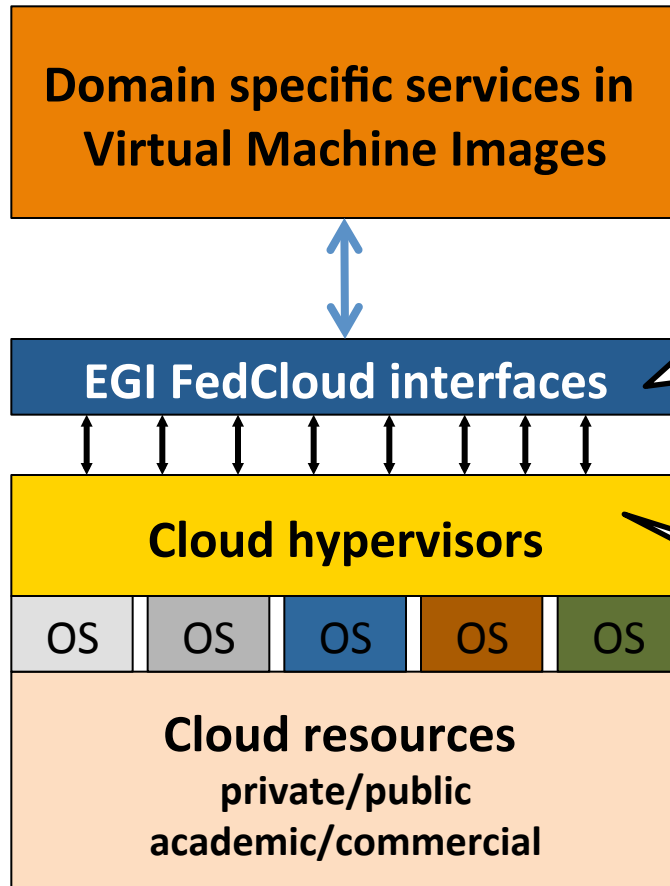Research
MArkers

DDN | GRIDScaler

- Distributed, federated storage and compute facilities
- Grid and Cloud compute platforms
- Virtual Research Environments
- > 200 user research projects

- 350 resource centres in 40 countries
- 400,000 logical CPU cores
- 190 PB disk, 180 PB tape
- > 99.6% reliability

**Domain specific services in Virtual Machine Images**

**EGI FedCloud interfaces**

**Cloud hypervisors**

OS | OS | OS | OS | OS

**Cloud resources**
private/public
academic/commercial
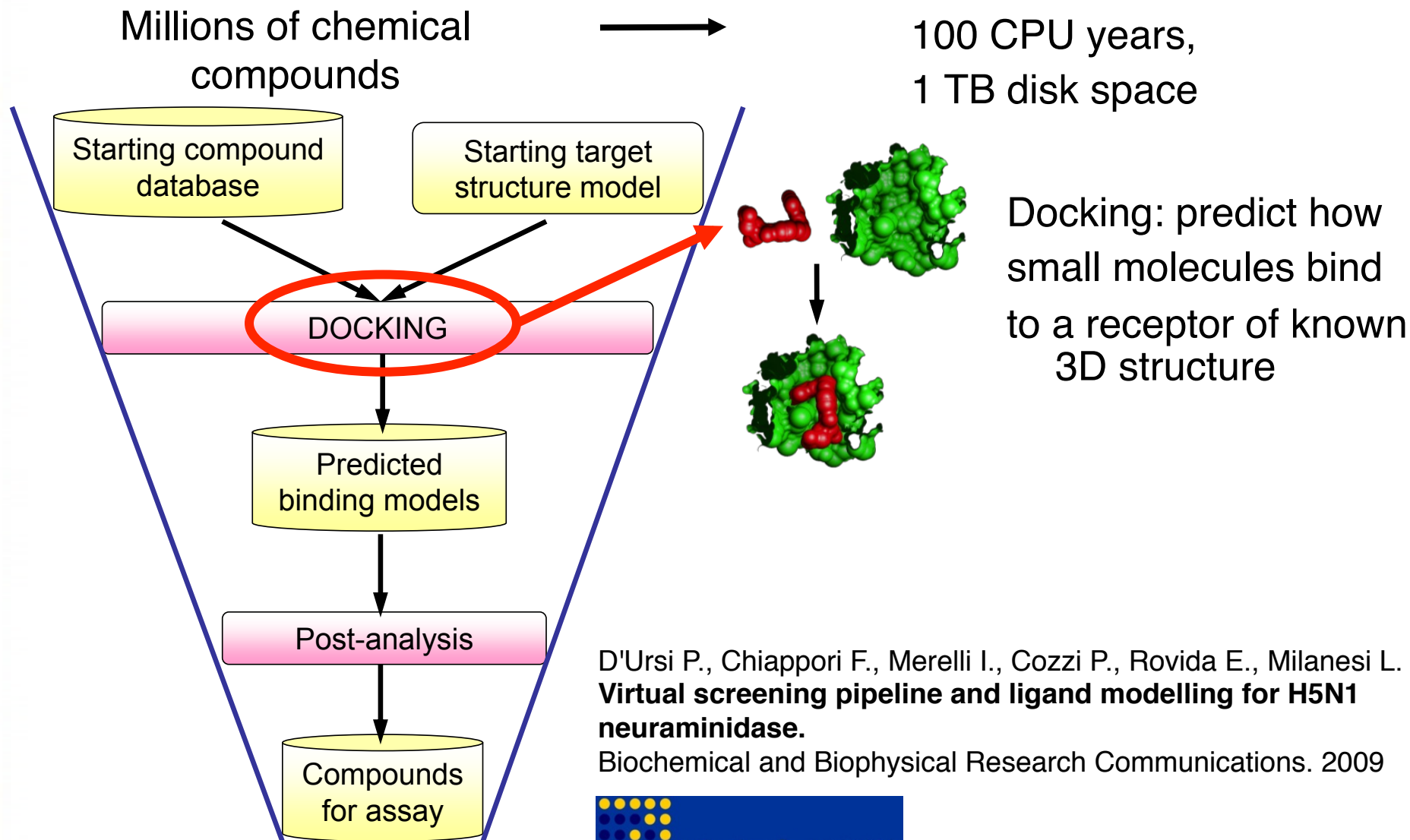
**http://go.egi.eu/cloud**

Standards enable federation
- OCCI: VM Image management
- OVF: VM Image format
- BDII: Information system
- X509: Authentication
- APEL: Accounting
- (CDMI: Cloud storage)
+ VM image Marketplace

Cloud hypervisor is a local choice. Eg.
- OpenStack
- OpenNebula
- EmotiveCloud (Spain)
- Okeanos (OpenStack impl. in GR)
- WNoDeS (Italy)
- …

Millions of chemical compounds →

100 CPU years,
1 TB disk space

Starting compound database

Starting target structure model

DOCKING

Docking: predict how small molecules bind to a receptor of known 3D structure

Predicted binding models

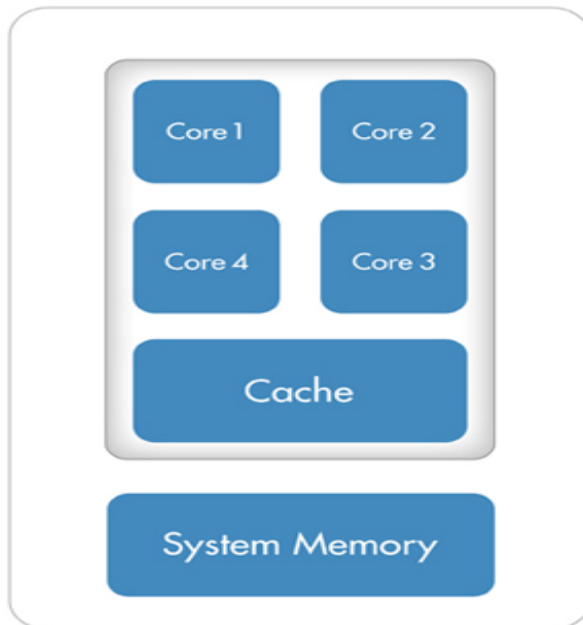Post-analysis

Compounds for assay

D'Ursi P., Chiappori F., Merelli I., Cozzi P., Rovida E., Milanesi L. **Virtual screening pipeline and ligand modelling for H5N1 neuraminidase.**
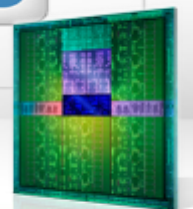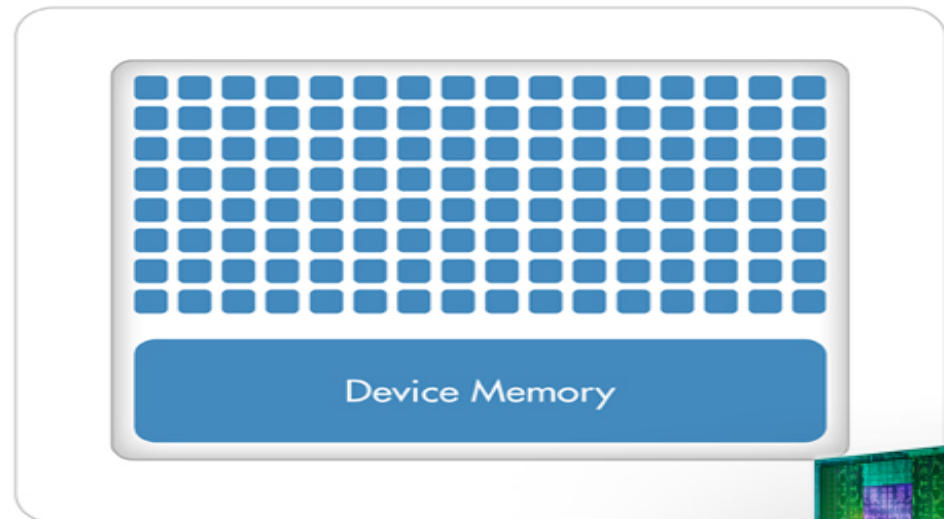Biochemical and Biophysical Research Communications. 2009

**BioinfoGRID**
Bioinformatics Grid Application for life science

GPUs implement a SIMD (Single Instruction Multiple Data) many-core architecture, providing a very high level of parallelism on intense data-parallel computation problems.

**CPU (Multiple Cores)**

| Core 1 | Core 2 |
|--------|--------|
| Core 4 | Core 3 |

Cache

System Memory

**GPU (Hundreds of Cores)**

Device Memory

- GPU-based solution in bioinformatics for:
  - Sequence Database Searching
    - CUDASW++
  - Multiple Sequence Alignment
    - CUDA-BLASTP
  - Next-Generation Sequencing
    - DecGPU, CUDA-EC, Musket, SOAP3-dp, CUSHAW
  - Genome-Wide Association Studies
    - Mendel_GPU, GENIE, SWIFTLINK
  - Motif Finding
    - mCUDA-MEME

# G-SNPM
## GPU SNP Mapping

- SNP genotyping analysis is very susceptible to SNPs chromosomal position errors;

- SNP mapping data are provided along the SNP arrays without information to assess in advance their accuracy;

- moreover, mapping data are related with a given build of a genome and need to be updated when a new build is available.

- The aim of **MIMOmics** is to develop new statistical methods for the integrated analysis for metabolomics, proteomics, glycomics and genomic datasets in large studies.

- Our partners are involvement involve in EU funded projects, i.e. **GEHA**, **IDEAL**, **Mark-Age**, **ENGAGE**, **EuroSpan**, **and BBMRI**



- In these consortia the primary goal is to **identify molecular profiles that monitor and explain complex traits** with novel findings so far.

- **MIMOmics** web site http://www.mimomics.eu at CNR (Milan, Italy)

MIMOmics resources
(data sets and computational tools)
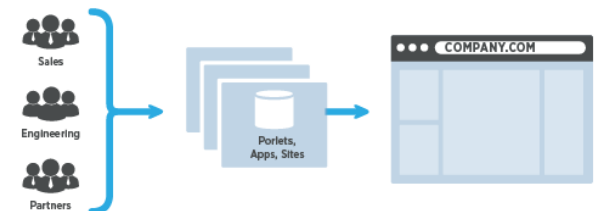
MIMOmics
authorized users



Project Web Portal to:
• create define the  users credentials for all  MIMOmics resources
• access MIMOmics resources
• develop, test and use tools on the data sets available
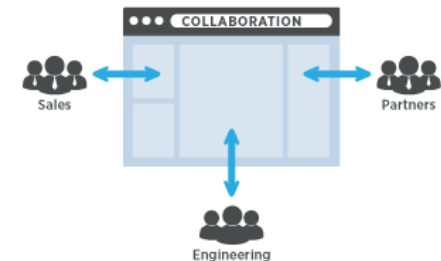• create pipeline of analysis combining tools and data sets

- The Omics Scientific Web Portal is based on **Liferay Portal tecnology**

- Liferay is a **robust** technology, fully supported in terms of **accessibility** and **scalability**

- Liferay provides a flexible template interface

- With Liferay the users can **manage contents** and documents in a distribuited and dinamic way over internet

- Liferay is compliabt with the **Java Portlet API 2.0**

**Documents Management**
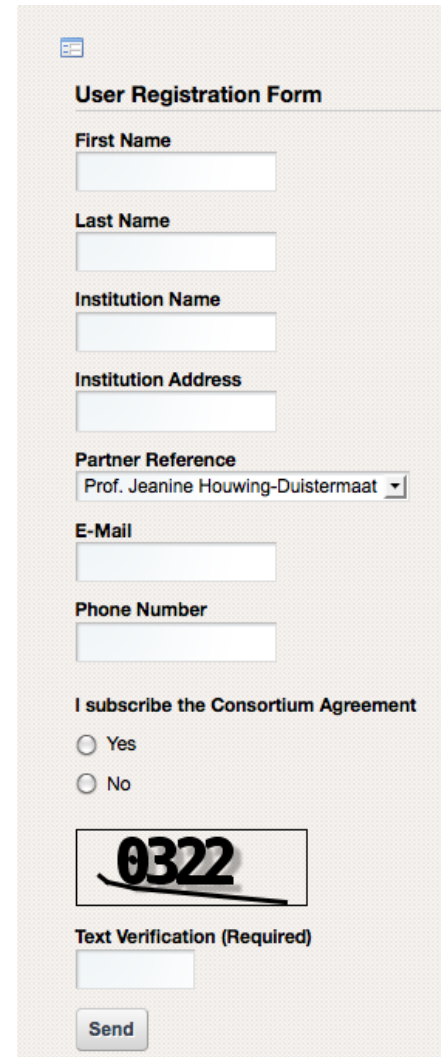
**Web Editing**
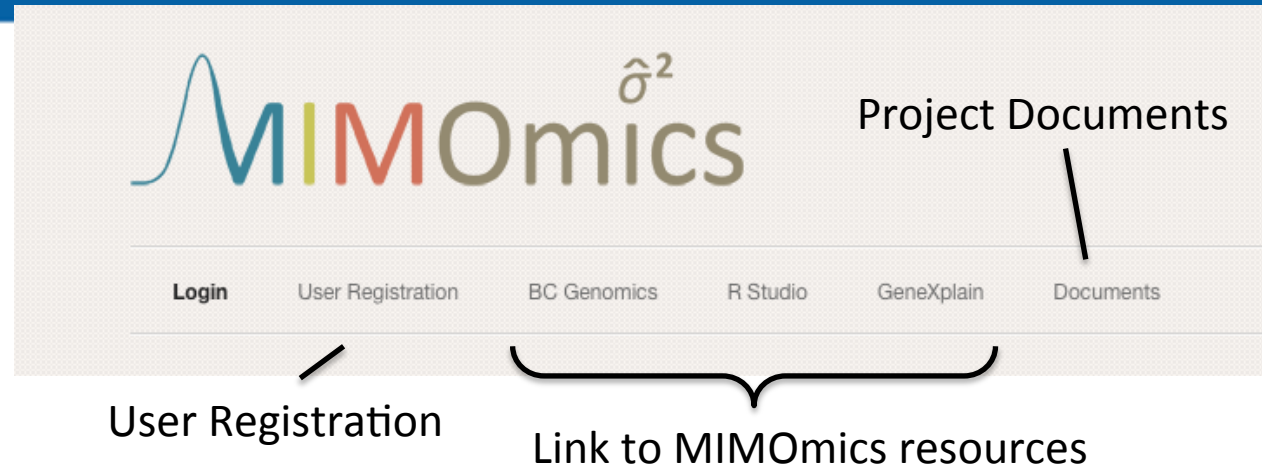
**Collaboration, Services**

User Registration

Project Documents

Link to MIMOmics resources
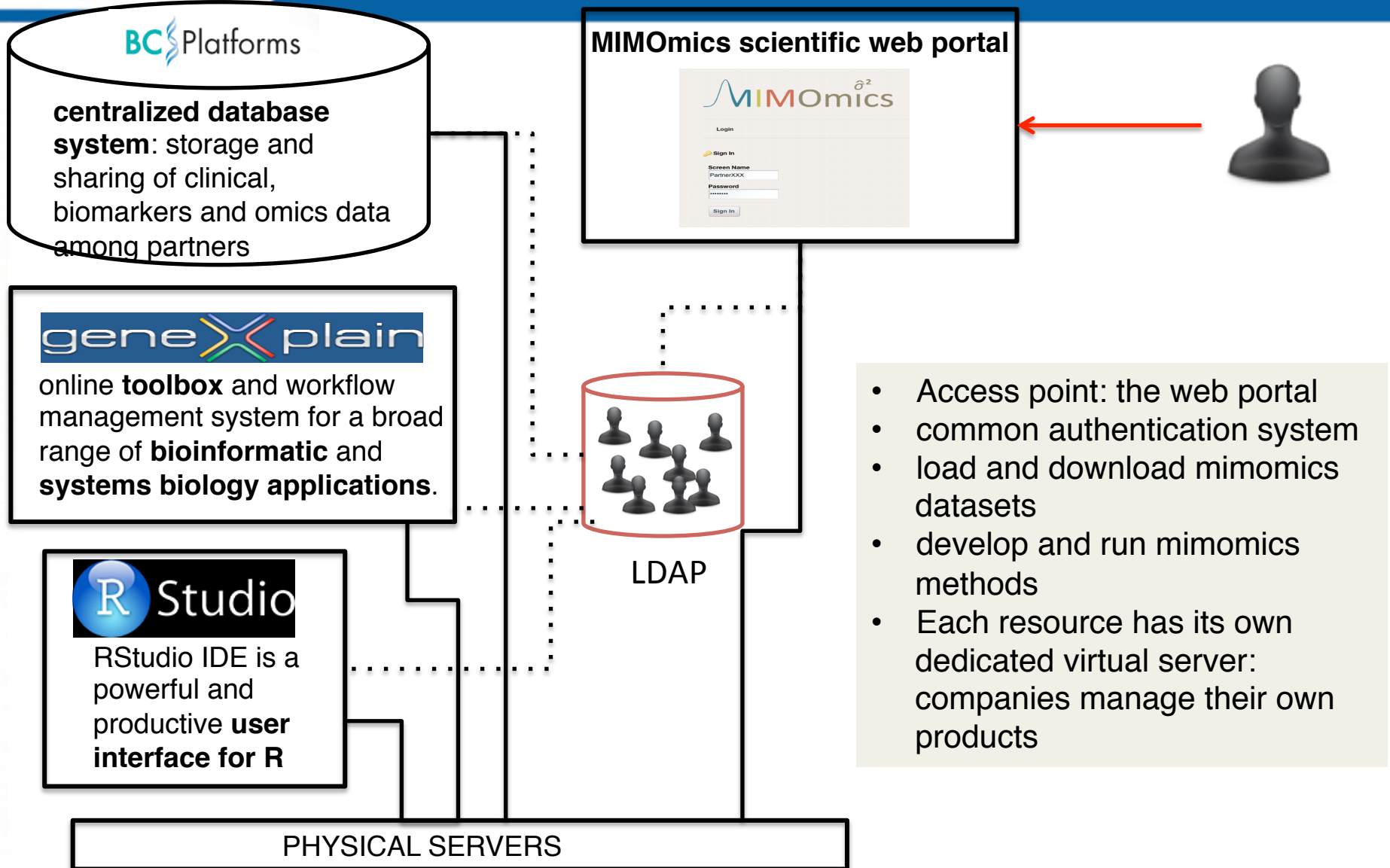
**Omics scientific web portal:**
- partner references can create new users with the same credentials for all MIMOmics resources
- access MIMOmics resources
- load and download MIMOmics datasets
- develop, test and use MIMOmics methods
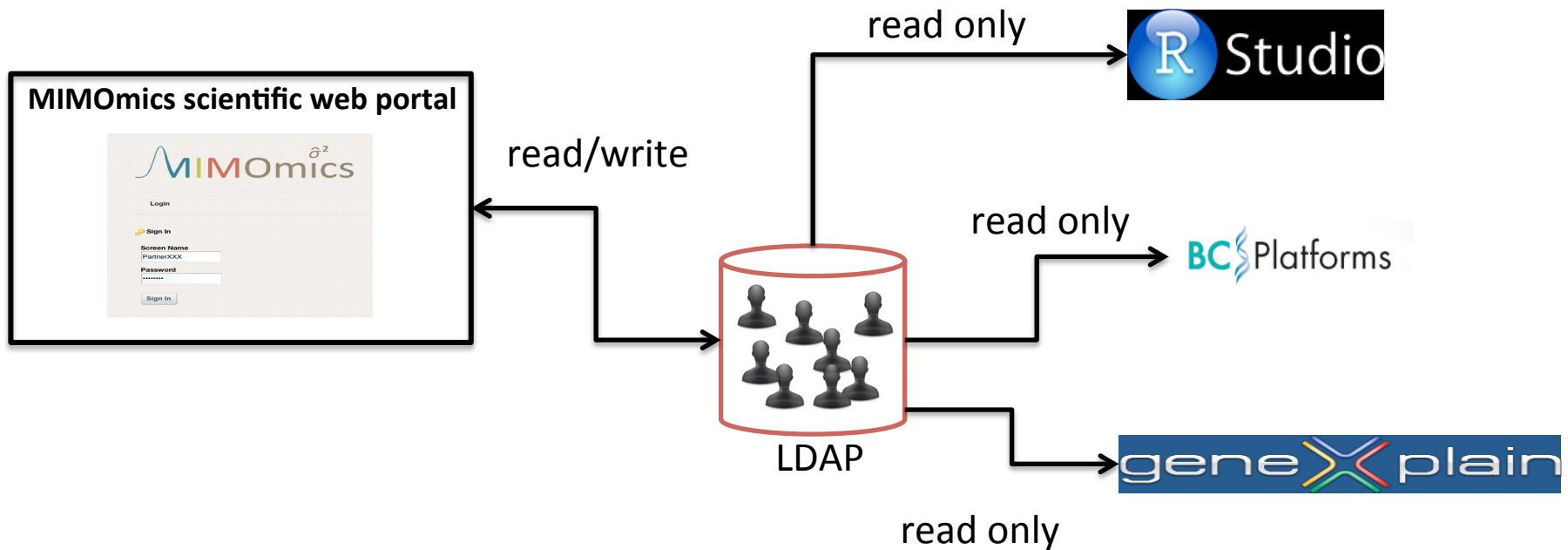- create pipeline of analysis combining tools and data sets

InterOmics
CNR

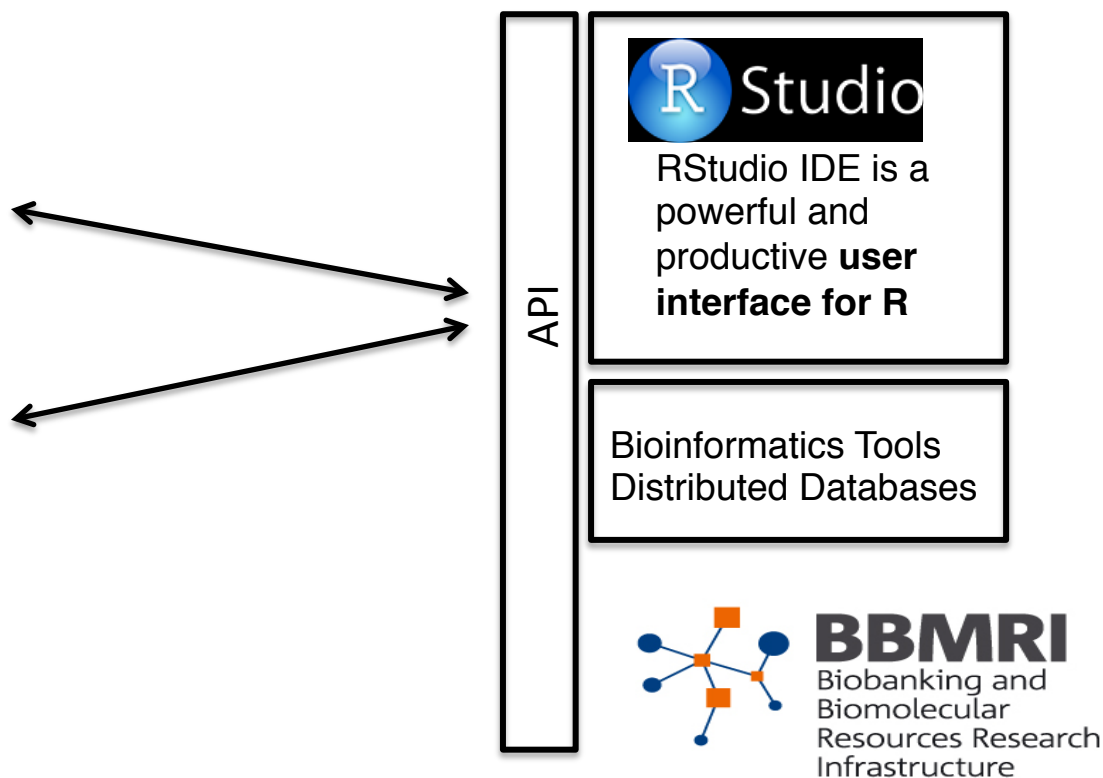**MIMOmics scientific web portal**

MIMOmics

Login

Sign In

Screen Name
PartnerXXX
Password

Sign In

**BC Platforms**

**centralized database system**: storage and sharing of clinical, biomarkers and omics data among partners

gene plain

online **toolbox** and workflow management system for a broad range of **bioinformatic** and **systems biology applications**.

R Studio

RStudio IDE is a powerful and productive **user interface for R**

LDAP

- Access point: the web portal
- common authentication system
- load and download mimomics datasets
- develop and run mimomics methods
- Each resource has its own dedicated virtual server: companies manage their own products
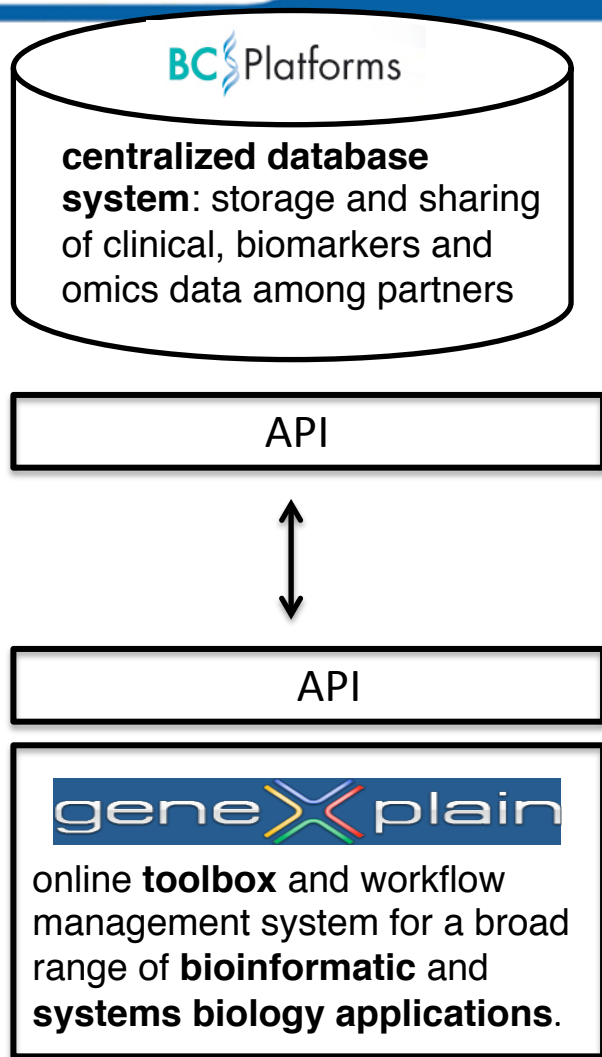
PHYSICAL SERVERS

InterOmics
CNR

- R packages available in RStudio server
  - core Bioconductor packages
  - R packages for multi-omics data analysis
    - **iCLuster**, a joint latent variable model for integrative clustering, (Shen et al., Bioinformatics, 2009)
    - **RISA**, converting experimental metadata from ISA-tab into Bioconductor data structures, (Gonzalez-Beltran et al., Bioconductor)
    - **OmicKriging**, Poly-Omic Prediction of Complex Traits, (Wheeler et al., 2013, arXiv:1303.1788)
    - *****ABEL,** facilitate statistical analyses of polymorphic genomes data (Yurii Aulchenko)
    - **iNEMO,** integration of NEtworks with Multi-Omics (E. Mosca, L. Milanesi)

Users are managed by the MIMOmics Scientific Web portal through the *Lightweight Directory Access Protocol* (LDAP).

**centralized database system**: storage and sharing of clinical, biomarkers and omics data among partners

Ad hoc API will be used for the integration of different resources in Cloud.

API

API

online **toolbox** and workflow management system for a broad range of **bioinformatic** and **systems biology applications**.

API

RStudio IDE is a powerful and productive **user interface for R**

Bioinformatics Tools Distributed Databases

**BBMRI**
Biobanking and Biomolecular Resources Research Infrastructure

Safebox set-up

**Several Omics Datasets:**
Genomics, Glycomics, Proteomics, Metabolomics/Lipidomics

**Several Studies**:
Aging, Cancer, Isolated Populations studies, Multiple Sclerosis, Obesity and Metabolic sSyndrome

**Biological Resource based on the BBMR standard Infrastructure:**

**BBMRI**
Biobanking and
Biomolecular
Resources Research
Infrastructure

- **SAM Tools** provide various utilities for manipulating alignments in the SAM format, including sorting, merging, indexing and generating alignments in a per-position format.

- **The Genome Analysis Toolkit** or **GATK** is a software package developed at the Broad Institute to analyse next-generation resequencing data.

- **Granvil**: Gene- or Region-based ANalysis of Variants of Intermediate and Low frequency

- **Annovar**: Functional annotation of genetic variants from high-throughput sequencing data.

- **PLINK** is a free, open-source whole genome association analysis toolset, designed to perform a range of basic, large-scale analyses in a computationally efficient manner.

- **IMPUTE** is a program for estimating ("imputing") unobserved genotypes in SNP association studies.

databases
- R Biomodels
- RW Biopath
- R Ensembl
- R EnsemblMouse
- R EnsemblRat
- R GO
- R GeneWays
- R HumanCyc
- R Reactome
- Utils

Load data

Normalize data

Detect differentially expressed genes

Discover functional enrichment

Identify master regulators in networks

Analyze regulatory genome regions

Analyze ChIP-Seq data

**Some of data and analysis tools based on GeneXplain**

# RStudio Server

- An instance of RStudio server has been installed and available for the MIMOmics users

- RStudio Integrated Development Environment is a powerful and productive user interface for R (http://www.rstudio.com/)

**Powerful productivity tools**

- Syntax highlighting, code completion, and smart indentation
- Execute R code directly from the source editor
- Easily manage multiple working directories using projects
- Quickly navigate code using typeahead search and go to definition

**An IDE built for R**

- Workspace browser and data viewer
- Plot history, zooming, and flexible image and PDF export
- Integrated R help and documentation
- Sweave authoring including one-click PDF preview
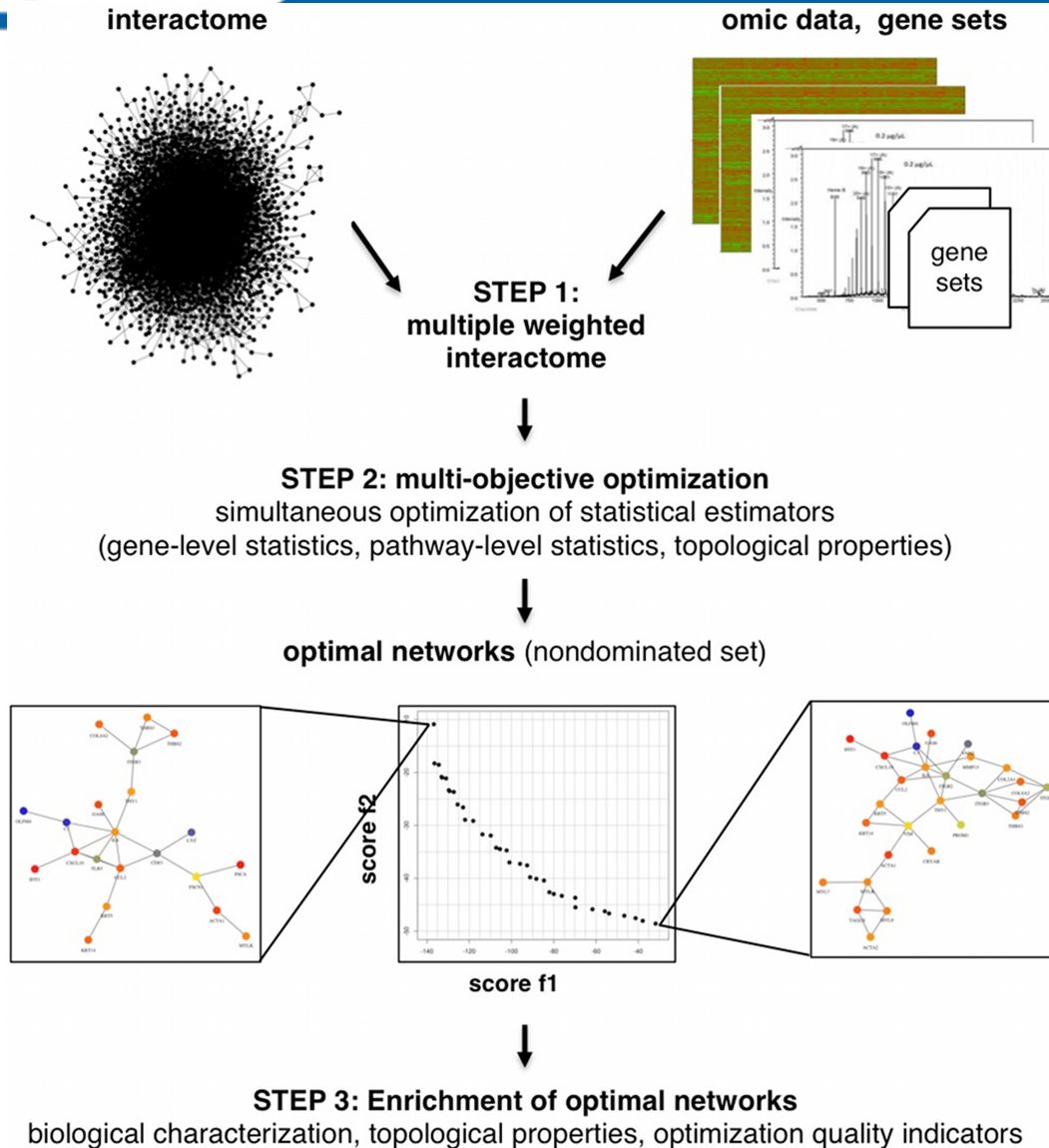- Searchable command history

# RStudio

## Examples of R packages for multi-omic data analysis:

- from the literature
  - **iCLuster**, a joint latent variable model for integrative clustering, (Shen et al., Bioinformatics, 2009)
  - **RISA**, converting experimental metadata from ISA-tab into Bioconductor data structures, (Gonzalez-Beltran et al., Bioconductor)
  - **OmicKriging**, Poly-Omic Prediction of Complex Traits, (Wheeler et al., 2013, arXiv:1303.1788)
  - **piano**, Platform for integrative analysis of omics data (Varemo, et al., 2013, NAR)

- from MIMOmics parters
  - *\***ABEL (GenABLE, OmicABLE, ProbABLE, … )** facilitate statistical analyses of polymorphic genomes data (Yurii Aulchenko)
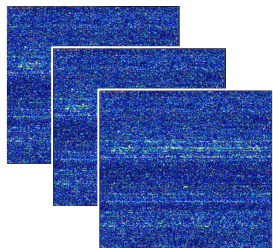  - **network-based integration of omics** (Mosca E, Milanesi L, *et al.* submitted)
  - Ecc.

**Integrating omic data**:
- Analyze the biological components and their interactions,
- Define a multiple-weighted **network**
- Find the **optimal modules** on the basis of the simultaneous optimization of **several statistical estimators**
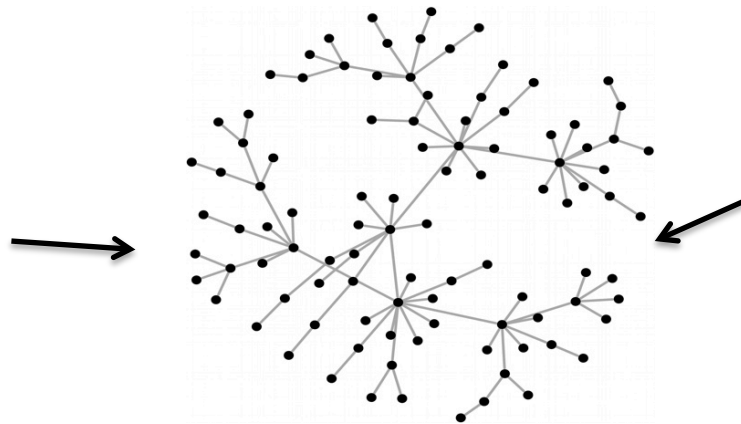
Mosca E, Milanesi L, *et al*.

**HCV and Host protein-protein interactions**

**Expression data of stepwise hepatocarcinogenic process**



GSE6764 (Geo Database)
Affymetrix HG-U133A
75 tissue samples
**Normal, Cirrhosis,
Dysplasia, Hepatocellular carcinoma**

### HCV - Host PPI

Kwofie SK et al. Infect Genet Evol 2011
DeChassey B et al. Mol Syst Biol. 2008
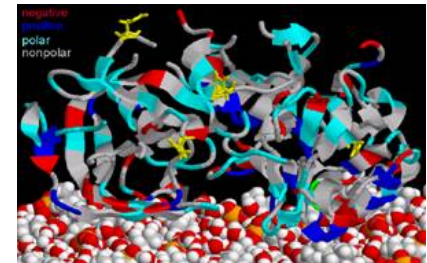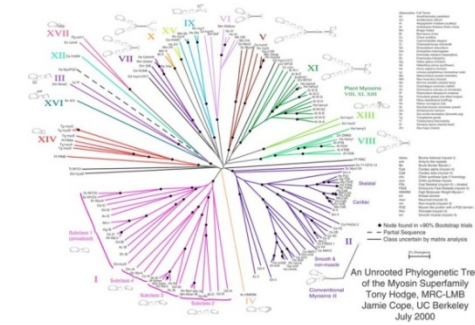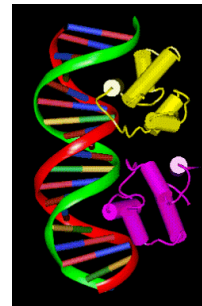
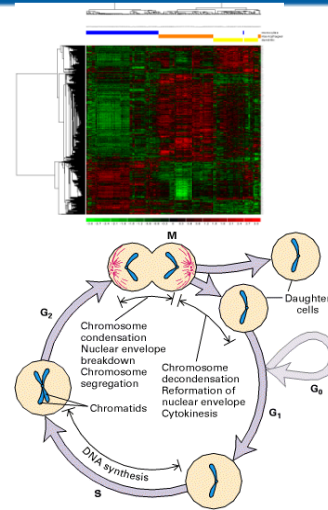Tot: 542 HCV- Host interactions

### Host PPI

Franceschini A et al. Nucl Acid Res 2013
STRING v9.1

Tot: ~224000 human interactions

**HCV – Host interactome
with multiple transcriptomic data**

**OBJECTIVE**
Identification of subnetworks enriched in differentially expressed
genes and HCV-host protein-protein interactions

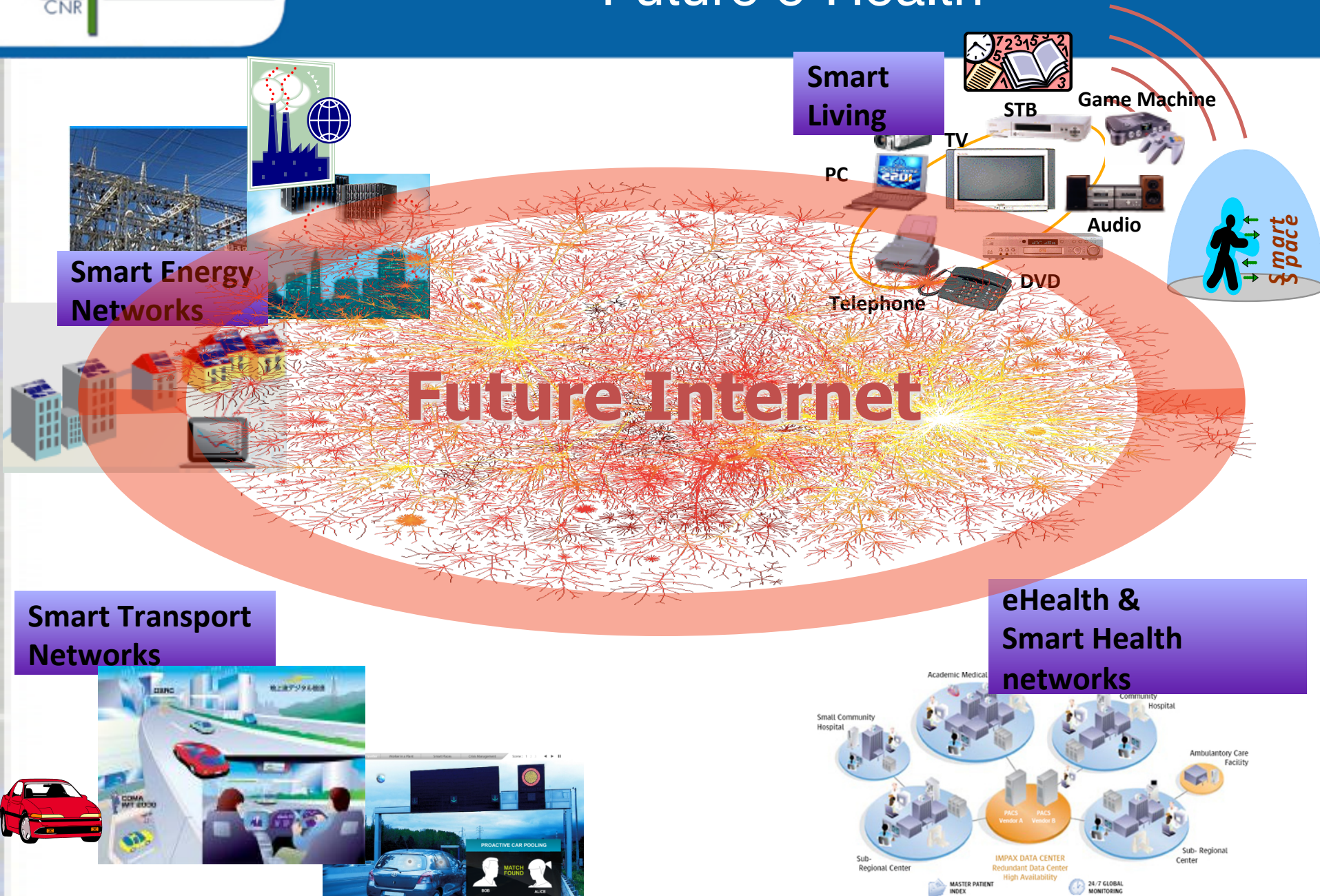- **Personalised medicine** will require sequencing of the genomes of large numbers of patients and volunteers

- It will be necessary to compare at least some of these genomes with the reference data collections

- Most hospitals and clinical research institutes will not wish to maintain up-to-date copies of the reference data collections

- It will be therefore be necessary to send these genomes to the institutes that hold the reference data collections

- It seems likely that this will be achieved using **secure VMs and secure clouds** holding the reference data collections

- EMBL-EBI is engaging with stakeholders to evaluate opportunities in this area.

- The use of **Big Data** and the **Omics technologies** will improve the research for the future **personalized system medicine** since the disease phenotypes arise from complex interactions among genetic factors and environment.

- The use of public's bioinformatics resources data center in connection with specialized **BioBanks** will be progressively used for **large-scale population biomarker discovery** and validation by integrating clinical and genetic databases and providing an integrated access to this huge amount of information.

- A range of new applications in biomedical data mining based on **Cloud Computing** are in fast development.

InterOmics
CNR

Ministero dell'Istruzion
dell'Università e Ricerca

EXPO
MILANO 2015

Local organizing committee:

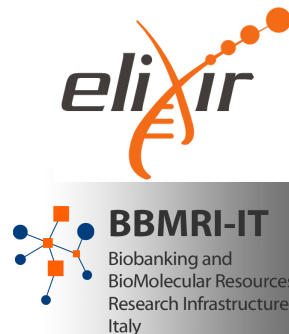M. Lavitrano, E. Bravo, MG Daidone, R. Lawlor, L. Milanesi, B. Parodi, D. Pistillo, G. Stanta.