# Efficient preprocessing and analysis of omics data:
## Experiences at University Magna Graecia of Catanzaro

**Mario Cannataro**

**Bioinformatics Laboratory &
Data Analytics Research Center,
Dep. of Medical and Surgical Sciences
University "Magna Græcia" of Catanzaro,
Italy**

**cannataro@unicz.it**

CGW 2017, Krakow, Poland – 23-25 October, 2017

# Brief History

- 1998: Format start of University of Catanzaro

- 2001/2002: Start of a Inter-University **Bachelor on Informatics and biomedical engineering**

- 2003: **University Campus** in its initial stage (right)

- 2004: Start of the interdisciplinary **PhD Program on Informatics and biomedical engineering**

- 2017: **University Campus** now (bottom)

**Core Facilities**
- **Genomics**
- **MicroArray**
- **Proteomics**
- **Mass Spectrometry**
- **Bioinformatics**
- **Nanotechnology**
- **Biomechatronics**

University of Catanzaro,
Campus "Salvatore Venuta",
Medicine and Bioscience Buildings

University Hospital

University Labs and Classrooms

Congress Hall

Master and PhD courses on Informatics and Biomedical Engineering,
- within the Faculty of Medicine
- collaboration with Tech Univ. of Milano, Univ. of Napoli "Federico II", Univ. of Calabria

Research groups,
- **Computer Engineering**
- **Control Systems & Electronic/Informatics Bio-engineering**
- **Mechanical Bio-Engineering**
- **Electronics**
- **Nanotechnology**

cannataro@unicz.it

CGW 2017, Krakow, 24 October 2017

- **Predictive, preventive, personalized and participatory (P4) medicine** is an emerging medical model that is based on the customization of all medical aspects (i.e. practices, drugs, decisions) of the individual patient
  - **Predictive medicine**: early diagnosis of a large set of diseases through regular scanning of the omic data (e.g. biomarker discovery)
  - **Preventive medicine**: prevent the development of some diseases before the appearance of symptoms through the monitoring of 'wellness', i.e. the behaviour of patients avoiding possible unhealthy practices
  - **Personalized medicine (Precision medicine)** is an emerging approach for disease treatment and prevention that takes into account individual variability in genes, environment, and lifestyle for each person. **Pharmacogenomics has a key role.**
    - January 20, 2015, President Obama announced the Precision Medicine Initiative® (PMI) - NIH
  - **Participatory medicine**: the active involvement of patients in the medical processes (e.g. patient-driven networks).
    - 'Participatory Medicine is a model of cooperative health care that seeks to achieve active involvement by patients, professionals, caregivers, and others across the continuum of care on all issues related to an individual's health.' (Society for Participatory Medicine)



P4 MEDICINE institute

Briefings in Bioinformatics Advance Access published September 8, 2015

Briefings in Bioinformatics, 2015, 1–9

doi: 10.1093/bib/bbv076
Paper

OXFORD

Methodologies and experimental platforms for generating and analysing microarray and mass spectrometry-based omics data to support P4 medicine

Pietro H. Guzzi, Giuseppe Agapito, Marianna Milano and Mario Cannataro
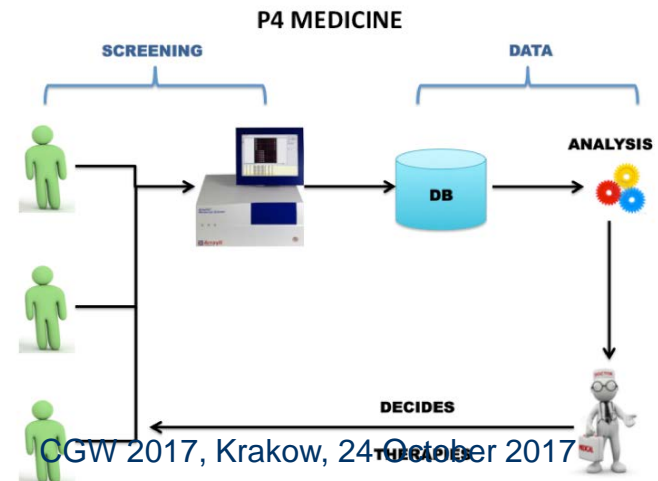
EXPERT REVIEW

Expert Review of Precision Medicine and Drug Development
Personalized medicine in drug development and clinical practice

ISSN: (Print) 2380-8993 (Online) Journal homepage: http://www.tandfonline.com/loi/tepm20

Experimental treatment of multiple myeloma in the era of precision medicine

Maria Teresa Di Martino, Mariamena Arbitrio, Pietro Hiram Guzzi, Mario Cannataro, Pierosandro Tagliaferri & Pierfrancesco Tassone

**P4 MEDICINE**

SCREENING | DATA | ANALYSIS | DECIDES

# OUTLINE

- Experiences at University of Catanzaro in high performance management, preprocessing and analysis of omics data

  - PART I: Genomics data

  - PART II: Proteomics data

  - PART III: Interactomics data

- Conclusions

*Dubium sapientiae initium*

UMG

*Dubium sapientiae initium*

- **Microarray-based Genomics Data**
  - gene expression and genotyping data
- **Automatic preprocessing of gene expression data**
  - **micro-CS**: preprocessing and annotation of gene expression data
- **Statistical analysis of genotyping data**
  - **DMET-Analyzer**: preprocessing and analysis of DMET (Drug Metabolizing Enzymes and Transporters) SNP (Single Nucleotide Polymorphism) data
  - **coreSNP** (multicore implementation)
- **Data Mining analysis of genotyping data**
  - **DMET-Miner**: Association Rule Mining to extract from DMET data Association Rules able to correlate the contemporary presence of SNPs with patient's conditions (e.g. TOX vs NOTOX)

- **Microarray-based Genomics Data**
  - gene expression and genotyping data
- Automatic preprocessing of gene expression data
  - **micro-CS**: preprocessing and annotation of gene expression data
- Statistical analysis of genotyping data
  - **DMET-Analyzer**: preprocessing and analysis of DMET (Drug Metabolizing Enzymes and Transporters) SNP (Single Nucleotide Polymorphism) data
  - **coreSNP** (multicore implementation)
  - **cloud4SNP** (cloud-based implementation)
- Data Mining analysis of genotyping data
  - **DMET-Miner**: Association Rule Mining to extract from DMET data Association Rules able to correlate the contemporary presence of SNPs with patient's conditions (e.g. TOX vs NOTOX)

*Dubium sapientiae initium*

# Gene Expression Microarray

– Sample Preparation
– Sample Deposition on Chip
– Ibridization
– Raw Data (Fluorescence Images)
– Preprocessing and Analysis



*MultiExperiment Viewer*

Analysis

WEKA
The University of Waikato

Samples

Sample annotation

Genes

Gene expression matrix

Gene annotation

Gene expression levels

**Diseased**
**Healthy**

Summarization, Normalization, Annotation

chip-specific libraries

Gene 1

Gene 2

CAMPIONE (donatore malato) (batterio patogeno) (pianta crioresistente)
RIFERIMENTO (donatore sano) (batterio normale) (pianta normale)

Marcatura del mRNA con coloranti fluorescenti

ROSSO    VERDE

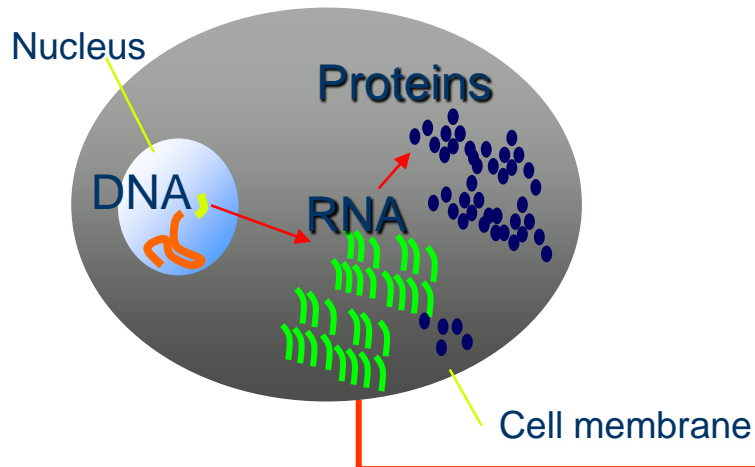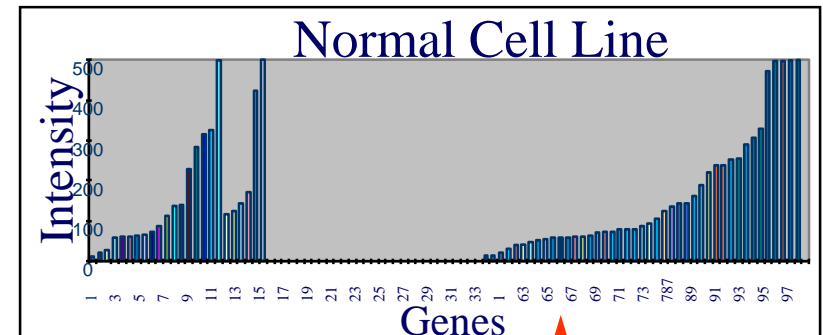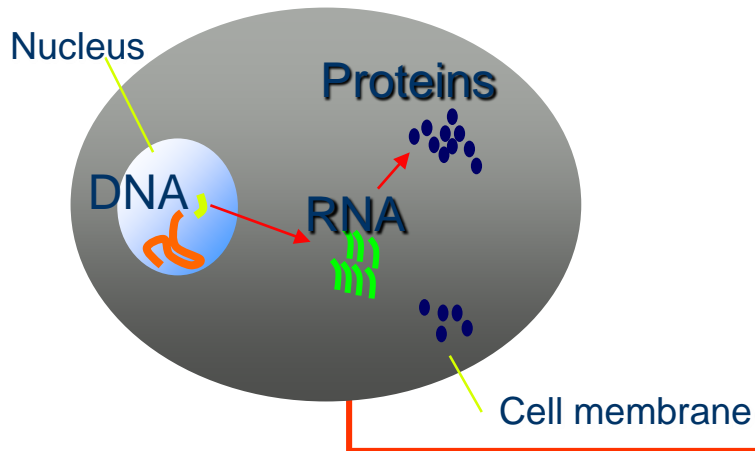Ibridizione tra il DNA marcato ed il DNA depositato sul chip
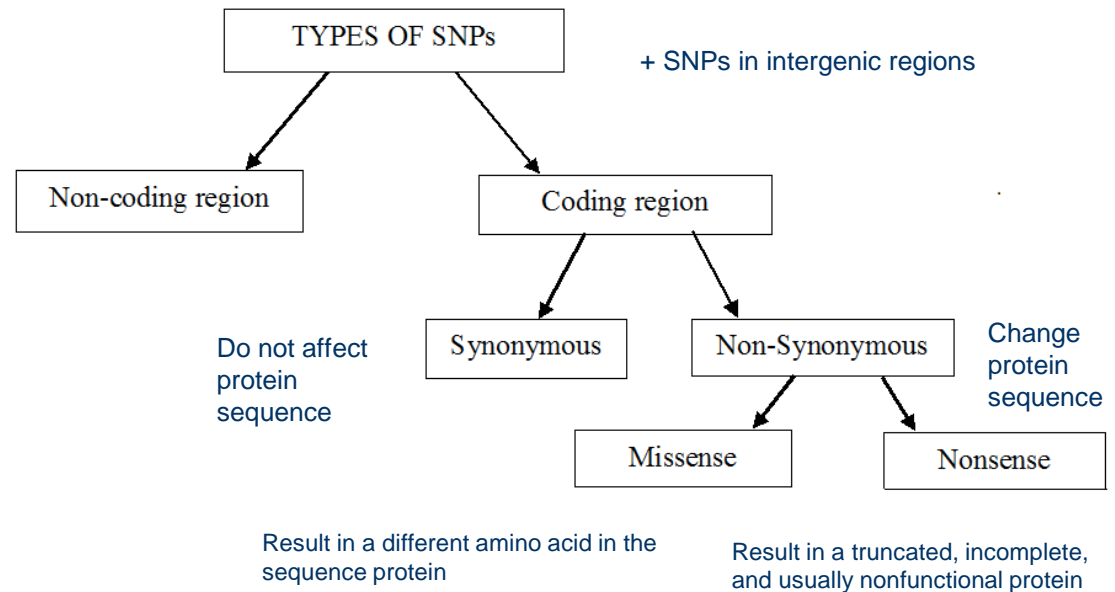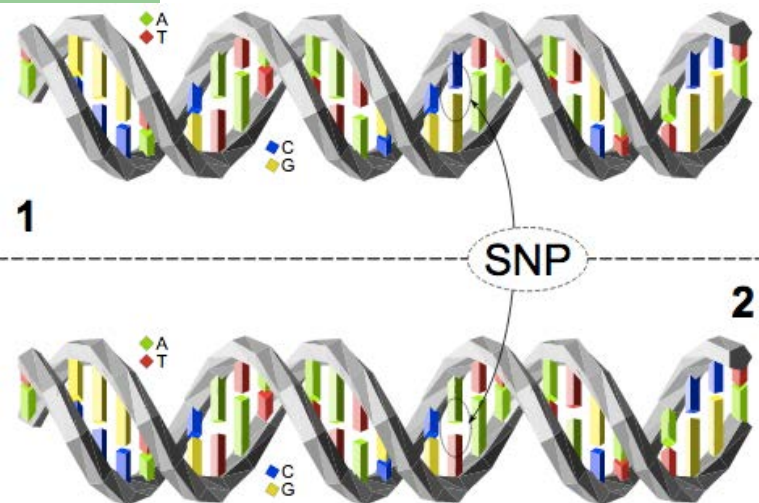
CHIP di DNA

cannataro@unicz.it

# Gene expression = RNA "volume"

# Single Nucleotide Polymorphisms (SNPs)

- A SNP is a variation in a single nucleotide that occurs at a specific position in the genome and is present in an appreciable percentage (e.g. >1%) within a population

- When in a specific base position, a base (e.g. C) appears in most individuals, but in a minority of individuals the position is occupied by a different base (e.g. A) then there is a SNP at that specific base position.



+ SNPs in intergenic regions

TYPES OF SNPs

Non-coding region — Coding region

Do not affect protein sequence

Synonymous — Non-Synonymous

Change protein sequence

Missense — Nonsense

Result in a different amino acid in the sequence protein

Result in a truncated, incomplete, and usually nonfunctional protein

cannataro@unicz.it

# Microarray-based Genomics data analysis

- We considered **gene expression** data and **SNP** (Single Nucleotide Polymorphism) **genotyping data** produced using microarrays

- We addressed two main problems:
  - The automation of the preprocessing pipeline of gene expression microarray data
    - **micro-CS**: preprocessing and annotation of gene expression data
  - The efficient (parallel) analysis of SNP genotyping microarray data
    - **DMET-Analyzer**: statistical analysis of DMET (Drug Metabolizing Enzymes and Transporters) SNP data
    - **coreSNP**: a parallel implementation of DMET-Analyzer
    - **DMET-Miner**: Association Rules correlating the presence of SNPs with patient's conditions (e.g. TOX vs NOTOX)

# Part I: Genomics Data

- Microarray-based Genomics Data

- Automatic preprocessing of gene expression data
    - **micro-CS**: preprocessing and annotation of gene expression data

- Statistical analysis of genotyping data
    - **DMET-Analyzer**: preprocessing and analysis of DMET (Drug Metabolizing Enzymes and Transporters) SNP (Single Nucleotide Polymorphism) data
    - **coreSNP** (multicore implementation)

- Data Mining analysis of genotyping data
    - **DMET-Miner**: Association Rule Mining to extract from DMET data Association Rules able to correlate the contemporary presence of SNPs with patient's conditions (e.g. TOX vs NOTOX)
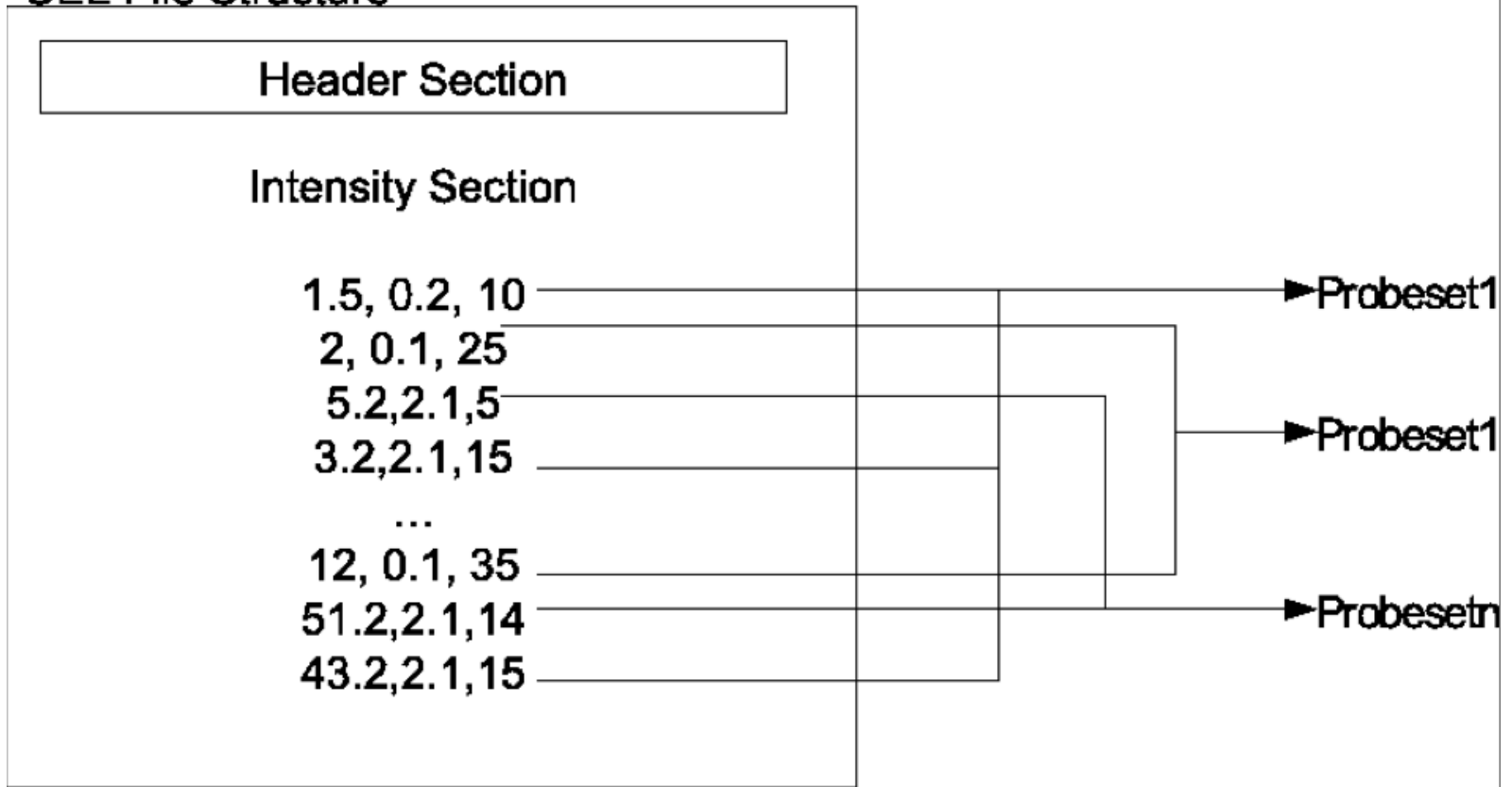
# Workflow of analysis of gene expression data

- **(i)  preprocessing**:
  - *Summarisation* recognizes the position of different genes in raw images, associating different regions of pixels to the unique gene that generated them.
  - *Normalisation* corrects the variation of gene expression in the same array due to experimental bias.
  - **Affymetrix Power Tools** (**APT**) is the command-line tool provided by Affymetrix for preprocessing. It uses some **chip-specific libraries**
- **(ii)  annotation**: associates each gene to functional information, such as biological processes, and a set of cross reference DB identifiers
  - Annotation is usually done **using vendor-provided libraries** (e.g. Affymetrix)
- **(iii)**  statistical or data mining **analysis**:
  - Analysis may be performed by external tools (e.g. TIGR MeV, Weka, R, etc.), but this requires data movement and further data reorganization
- **(iv)**  biological **interpretation** and access to **external knowledge bases**

**OPEN PROBLEM: THIS PROCESS INVOLVES SEVERAL TOOLS AND CHIP-SPECIFIC / VENDOR-SPECIFIC LIBRARIES AND MAY BE ERROR PRONE**

# The format or raw .cel files

CEL File Structure

Header Section

Intensity Section

1.5, 0.2, 10 → Probeset1
2, 0.1, 25
5.2,2.1,5
3.2,2.1,15 → Probeset1
…
12, 0.1, 35
51.2,2.1,14 → Probesetn
43.2,2.1,15

# Affymetrix Power Tools

- The preprocessing of Affymetrix CEL files is performed using specialized tools, e.g. the APT (Affymetrix Power Tools) command line tools

- **apt-probeset-summarize** is an APT tool for summarizing and normalizing expression probe data from CEL files.
  - It needs library files (either a cdf file or pgf/clf files) for defining probesets
  - It can perform different types of summarization, e.g. the Robust Multiarray Average (RMA) and the Probe Logarithmic Intensity Error (PLIER) algorithms
    - For exon arrays, the Detection Above BackGround (DABG) is provided
  - It can perform two main types of normalization, the quantile algorithm and the sketch-quantile algorithm, that re-scale each expression value of a dataset.

    ```
    apt-probeset-summarize -a rma-sketch -a plier-mm-sketch -d chip.cdf
    -o output-dir --cel-files cel_list.txt
    ```

    ```
    apt-probeset-summarize -a rma -d HuEx-1_0-st-v2.cdf -o/home/output -
    cel-files/home/list.txt
    ```

- A further step is **annotation,** that allows to associate to each expression value the related gene and further biological annotation,
  - e.g. database identifier, description of molecular function, associated protein domains, Gene Ontology data,

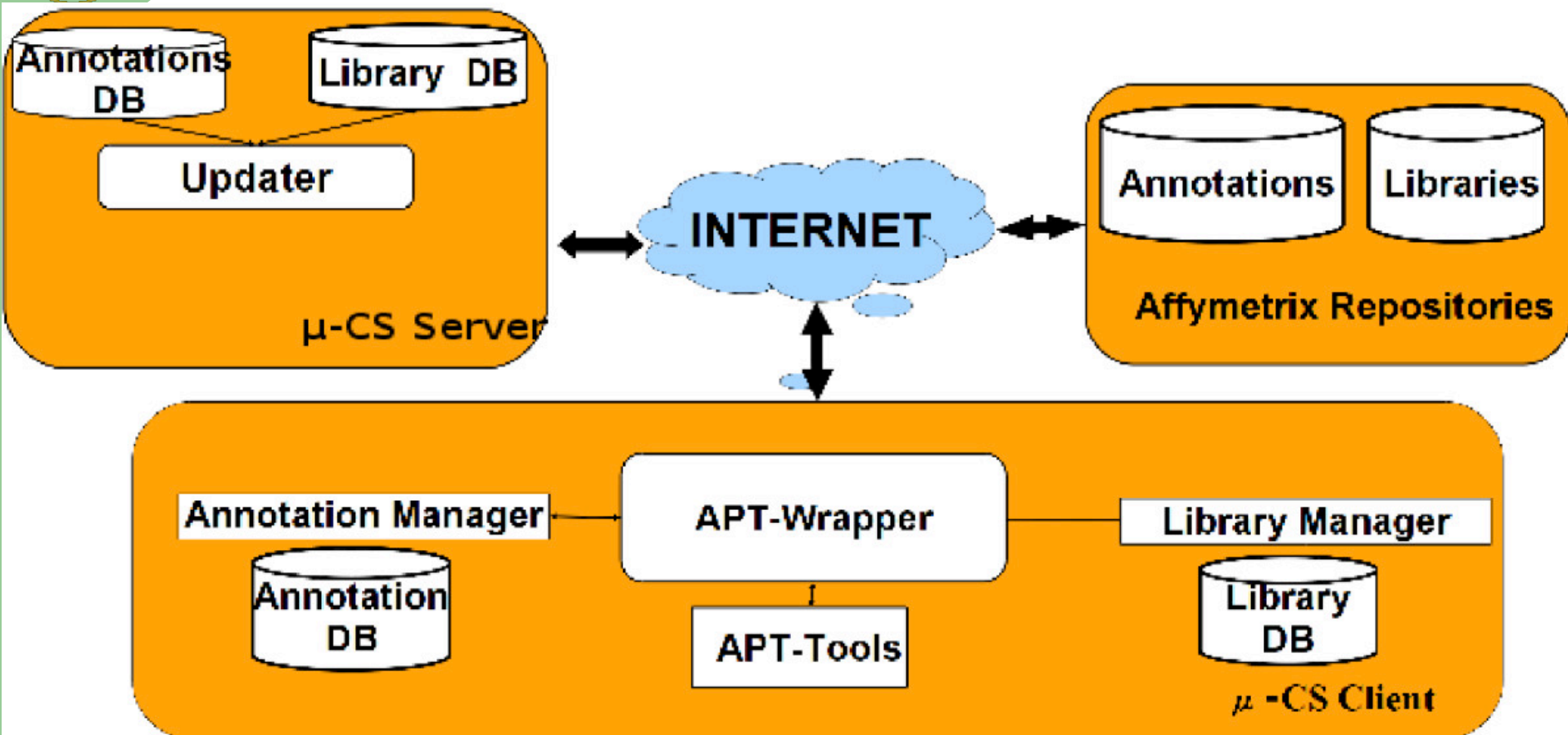- by using annotation files often provided by chip manufacturer

- The main drawbacks of such an approach are the need:
  - to generate and store intermediate files in a manual way that prevent the automation of the process;
  - to know details of the used chips and related preprocessing tools and libraries;
  - to manually download the most updated summarization and annotation libraries from the vendor website, and
  - to manually import preprocessed files in the analysis platforms (e.g. TM4 MeV, Weka, etc.), that may introduce errors.
- **The automation of the preprocessing pipeline could speed-up the entire analysis process and reduce possible errors, allowing the user to concentrate on biological aspects.**
- **micro-CS** is a tool to normalize, summarize and annotate gene expression data produced by Affymetrix microarray reducing user intervention
  - system handles the guided selection and automatic upgrading of the software libraries needed by Affymetrix tools to preprocess and annotate gene expressions
  - the tool is available with a stand-alone user-interface or as a TM4 MeV plug-in
    - TM4: a free, open-source system for microarray data management and analysis. Biotechniques. 2003 Feb;34(2):374-8.

- $\mu$-CS preprocesses Affymetrix microarray data to automatize summarization, normalization, and annotation of microarray.
  - $\mu$-CS users can directly manage binary data without worrying about locating and invoking the proper preprocessing tools and chip-specific libraries.
  - users of the $\mu$-CS plugin for TM4 can manage Affymetrix binary files without using external tools, such as APT (Affymetrix Power Tools) and related libraries.
- $\mu$-CS offers four main advantages:
  - (i) it avoids to waste time for searching the correct libraries,
  - (ii) it reduces possible errors in the preprocessing and further analysis phases, e.g. due to the incorrect choice of parameters or the use of old libraries,
  - (iii) it implements the annotation of preprocessed data,
  - (iv) it may enhance the quality of further analysis since it provides the most updated annotation libraries.
- $\mu$-CS client is freely available (TM4 plugin or standalone tool)

# micro-CS Architecture

Guzzi and Cannataro *BMC Bioinformatics* 2010, **11**:315
http://www.biomedcentral.com/1471-2105/11/315

BMC Bioinformatics

Open Access

**SOFTWARE**

$\mu$-CS: An extension of the TM4 platform to manage Affymetrix binary data

Pietro H Guzzi*[1] and Mario Cannataro*[1,2]

# Update of the Client Databases



References to Annotations DB
References Library DB

Updater

µ-CS Server

2) The µ-CS web server sends to the client most updates links to libraries and annotations.

Annotation Manager
Annotation DB
APT-Wrapper

Library Manager
Library DB
APT-Tools

µ-CS Client

1) The µ-CS tool requests the most updated annotation and library files to the µ-CS web server.



AnnotationDB    LibraryDB

Updater

µ-CS Server

1) Check for most updated summarization/ annotations libraries.

Affymetrix Repositories
www.affymetrix.com

Annotations    Libraries

2) Download of references to the newest summarization/annotation libraries

# Update of the Server Databases

(1) User can load transparently Cel files directly from the TMeV Loader (Madam).

(2) Then he/she has to set main parameters of summarization (e.g. kind of summarization and related parameters)

(3) Finally, processed files are ready to be exported in a file or loaded in TMeV

**Differential transcriptional response to cisplatinum in BRCA1-defective *versus* BRCA1-reconstituted breast cancer cells by microarrays.**

MT Di Martino[1,2], M Ventura[1], PH Guzzi[3], A Pietragalla[1], P Neri[2], A Bulotta[1], T Calimeri[1], V Barbieri[1,2], M Caraglia[4], P Veltri[3], M Cannataro[3], P Tassone[1,2], and P Tagliaferri[1,2]

## M19  WHOLE GENE EXPRESSION PROFILING SHOWS A DIFFERENTIAL TRANSCRIPTIONAL RESPONSE TO CISPLATINUM IN BRCA-1 DEFECTIVE VERSUS BRCA1-RECONSTITUTED BREAST CANCER CELLS

Di Martino MT[1,2], Guzzi PH[3], Ventura M[1], Pietragalla A[1], Neri P[2], Bulotta A[1], Calimeri T[1], Barbieri V[1,2], Caraglia M[4], Veltri P[3], Cannataro M[3], Tassone P[1,2] and Tagliaferri P[1,2]

FOCUS

## Automatic summarisation and annotation of microarray data

Pietro H. Guzzi · Maria Teresa Di Martino · Giuseppe Tradigo · Pierangelo Veltri · Pierfrancesco Tassone · Pierosandro Tagliaferri · Mario Cannataro

- To study the molecular bases of BRCA1-related differential sensitivity to the drug, we analyzed the whole gene expression profile of HCC1937 and HCC1937/wtBRCA1 cells following in vivo and in vitro exposure of tumor cells to cisplatinum (CDDP)
  - HCC1937 is a BRCA1-defective breast cancer cell line which discloses higher sensitivity to CDDP as compared to the BRCA1 full length cDNA transfected clone HCC1937/wtBRCA1
  - Gene expression profiling was performed by Affymetrix technology using **Human GeneArray 1.0ST**.
  - **Array data were preprocessed using μ-CS** and analyzed using Gene Expression Console, GeneSpring and Ingenuity Pathway Analysis (IPA).

# Part I: Genomics Data

- Microarray-based Genomics Data

- Automatic preprocessing of gene expression data
  - **micro-CS**: preprocessing and annotation of gene expression data

- Statistical analysis of genotyping data
  - **DMET-Analyzer**: preprocessing and analysis of DMET (Drug Metabolizing Enzymes and Transporters) SNP (Single Nucleotide Polymorphism) data
  - **coreSNP** (multicore implementation)

- Data Mining analysis of genotyping data
  - **DMET-Miner**: Association Rule Mining to extract from DMET data Association Rules able to correlate the contemporary presence of SNPs with patient's conditions (e.g. TOX vs NOTOX)

**UMG**

*Dubium sapientiae initium*

- The **Affymetrix DMET** (Drug Metabolism Enzymes and Transporters) platform **investigates 225 ADME genes**, i.e. genes involved in Absorption, Distribution, Metabolism and Excretion (ADME) of drugs, **for pharmacogenomics case-control study**.

- **Pharmacogenomics** is a branch of genomics that aims to predict the response to drugs of an individual based on an his/her genotype

  – **The hypothesis of DMET analysis is that a different response to drugs may be related to modifications (SNPs) in those ADME genes**

- We want identify the SNPs that make effective/ineffective a drug treatment using efficient Statistical Analysis and Association Rule Mining

| A | B | C | D | E | F | G | H | I | |
|---|---|---|---|---|---|---|---|---|---|
| PROBEID | NONRESP | NONRESP | NONRESP | NONRESP | NONRESP | NONRESP | NONRESP | NONRESP | N |
| AM_10001 | C/C | C/C | C/C | C/C | C/C | C/C | C/C | C/C | C |
| AM_10002 | G/G | G/G | G/G | G/G | G/G | G/G | G/G | G/G | |
| AM_10003 | T/T | C/T | C/T | T/T | C/T | C/T | C/C | C/C | C |
| AM_10004 | G/G | G/G | G/G | G/G | G/G | G/G | G/G | G/G | |
| AM_10005 | T/T | C/T | C/T | T/T | C/T | C/T | C/C | C/C | |
| AM_10006 | T/T | C/T | C/T | T/T | C/T | C/T | C/C | C/C | C |
| AM_10008 | A/A | A/A | A/A | A/A | A/A | A/A | A/A | A/A | A |
| AM_10010 | C/C | C/T | C/T | C/C | C/T | C/T | T/T | T/T | T |
| AM_10011 | A/A | A/G | A/G | A/A | A/G | A/G | G/G | G/G | |
| AM_10012 | G/G | G/G | G/G | G/G | G/G | G/G | G/G | G/G | |
| AM_10013 | A/A | A/G | A/G | A/A | A/G | A/G | G/G | G/G | |
| AM_10014 | C/C | C/T | C/T | C/C | C/T | C/T | T/T | T/T | |
| AM_10016 | T/T | C/T | C/T | T/T | C/T | C/T | C/C | C/C | |
| AM_10017 | G/G | G/G | G/G | G/G | G/G | G/G | G/G | G/G | |
| AM_10019 | A/A | A/A | A/A | A/A | A/A | A/A | A/A | A/A | A |
| AM_10020 | C/C | C/T | T/T | C/C | C/C | C/T | T/T | C/C | |
| AM_10021 | A/A | A/A | A/A | A/A | A/A | A/A | A/A | A/A | |
| AM_10022 | G/G | G/G | G/G | G/G | G/G | G/G | G/G | G/G | |
| AM_10023 | A/A | A/A | A/A | A/A | A/A | A/A | A/A | A/T | |
| AM_10024 | A/A | A/G | G/G | A/A | A/G | G/G | G/G | A/A | |
| AM_10025 | A/A | A/G | G/G | A/A | G/G | G/G | G/G | G/G | |
| AM_10028 | C/C | C/C | C/C | C/C | C/C | C/C | C/C | C/C | |
| AM_10030 | C/C | C/C | C/C | C/C | C/C | C/C | C/C | C/C | |
| AM_10031 | T/T | G/G | T/T | G/G | G/G | G/T | G/G | G/T | |
| AM_10033 | G/G | G/G | G/G | A/A | A/G | G/G | G/G | G/G | |
| AM_10034 | C/C | C/C | C/C | C/C | C/C | C/C | C/C | C/C | |
| AM_10035 | T/T | G/G | T/T | G/G | G/G | G/T | G/T | T/T | C |
| AM_10039 | T/T | T/T | T/T | T/T | T/T | T/T | T/T | T/T | T |
| AM_10042 | -/- | -/- | -/- | -/- | -/- | -/- | -/- | -/- | |
| AM_10044 | A/A | A/A | A/A | A/A | A/A | A/A | A/A | A/A | A |
| AM_10047 | C/C | C/C | C/C | C/C | C/C | C/T | C/C | C/T | |

# Workflow of a DMET pharmacogenomics experiment



| Probes | $Subject_1$ | $Subject_2$ | $Subject_3$ | $Subject_4$ |
|---|---|---|---|---|
| $Probe_1$ | C/C | C/C | T/T | T/T |
| $Probe_2$ | G/C | C/C | -/T | A/A |
| $Probe_3$ | C/T | C/T | C/T | C/T |
| $Probe_n$ | G/G | A/G | G/G | G/G |

- **OPEN PROBLEMS:**
    - **NO BIOINFORMATICS SOFTWARE FOR ANALYZING DMET SNP DATA WAS AVAILABLE AT THE TIME OF RESEARCH (ONLY MANUAL ANALYSIS WITH EXTERNAL SOFTWARE)**
    - **PERFORMANCE ISSUES WHEN INCREASING SIZE OF EXAMINED POPULATION**

# Workflow of a DMET pharmacogenomics experiment

- **Sample collection and DMET chip preparation**: biological samples are collected and treated to perform microarray experiments;
  - Affymetrix DMET chip allows the investigation of 1936 probes, each one representing a portion of the genome having a role in drug metabolism;
  - usually samples are divided in two classes, e.g. on the basis of the response to a drug (Population A and B).
- **DMET microarray experiments**: produce raw expression data, contained in .CEL files (one file per sample);
- **DMET data preprocessing and SNPs detection**: Affymetrix proprietary tools, (e.g. **apt-dmet-genotype** command line tool or **DMET Console**) are used to summarize and average expression values contained in .CEL file to produce .CHP files (one file per sample);
  - all .CHP files need to be combined to form a single table containing the detected SNPs for all samples;
  - similarly to gene expression microarray, the SNPs table will contain in position (i,j) the SNP detected by probe i in the sample j;
- **SNPs analysis**: usually the Fisher statistical test is used to evaluate if a different distribution of SNPs among two classes of samples is statistically significant or not. The analysis, often performed manually using external software, produces a list of SNPs ordered by p-value

**U M G**

*Dubium sapientiae initium*

- DMET-Analyzer is a tool for automated pre-processing and statistical analysis of SNP data generated by the Affymetrix DMET platform
  - the system finds and highlights all statistically significant SNP variations found in the data
- DMET-Analyzer supports the automatic statistical analysis in DMET-based pharmacogenomics studies.
  - It has a simple graphical user interface that allows users (doctors/biologists) to upload and analyze DMET files produced by DMET Console in an interactive way.
- Starting from a DMET dataset, DMET-Analyzer is able to find statistically relevant subsets of SNPs that separate two input classes
  - DMET-Analyzer tests in an automatic way if the distributions of each SNPs detected on each probe on the two classes of subjects (e.g. Healthy vs Diseased, RESP vs NORESP, TOX vs NOTOX) are statistically significant.
- The system is freely available and currently used by the Oncology Unit of the "Mater Domini" University Hospital, Catanzaro, Italy.

# Algorithm

DMET-Analyzer allows to test in an automatic way if the different distributions of each SNPs detected on each probe on the two classes of subjects (e.g. RESP vs NORESP) are statistically significant.

- It performs a massive amount of Fisher's tests without user intervention
- It provides several statistical corrections (FDR, Bonferroni) useful for low samples experiments
- It uses two parameters to
  - discard probes whose SNP distributions are "similar", avoiding useless tests
  - discard Fisher's tests with high p-value

| Probes | $Subject_1$ | $Subject_2$ | $Subject_3$ | $Subject_4$ |
|--------|-------------|-------------|-------------|-------------|
| $Probe_1$ | C/C | C/C | T/T | T/T |
| $Probe_2$ | G/C | C/C | -/T | A/A |
| $Probe_3$ | C/T | C/T | C/T | C/T |
| $Probe_n$ | G/G | A/G | G/G | G/G |



Fig. 2. DMET-Analyser flowchart.

cannataro@unicz.it

CGW 2017, Krakow, 24 October 2017

*Dubium sapientiae initium*

**a** User chooses the Data File to be loaded.

**b** User selects the classes for each sample by clicking on data table.

**c** Allele frequencies for each probe are calculated and shown to the user

Results are calculated

**d** User selects the Statistical Test and may visualize distribution.

**e** Results are shown to the user. He may visualize annotations and links to esternal databases.

# DMET-Analyzer: automatic analysis of Affymetrix DMET Data

BMC Bioinformatics

Pietro Hiram Guzzi[1*], Giuseppe Agapito[1], Maria Teresa Di Martino[2], Mariamena Arbitrio[3], Pierfrancesco Tassone[2], Pierosandro Tagliaferri[2] and Mario Cannataro[1]

*Dubium s*

# DMET-Analyzer use: loading input data and selection of classes



cannataro@unicz.it   Start Preprocess

# Heat Map probe navigator window



**Frequecy Differences Table - differences of SNPs frequencies for each probe (class B w.r.t class A)**



cannataro@unicz.it

*Dubium sapientiae initium*

# Statistical tests:



1) Setting correction;

2) Executing Fisher's tests;

3) Analyze results.

- The Fishers test calculator uses the data contained into two Occurrences Tables to compute the Fishers test for each couple of allele belonging to the two classes.
- The algorithm avoids the computation of trivial tests, e.g. tests where the Fisher Test contingency table contains three zero values that leads to a p_value=1 are discarded.
- Fisher tests results with *p_value > Ft* (**Filter Threshold**, default 0.05) are discarded.



cannataro@unicz.it

## TABLE 1
### Example Dataset Containing Alleles Detected in 8 Subjects Through 2 Probes ($np = 2$ and $ns = 8$)

**SNPs Input Table**

| Probes | $S_1$ | $S_2$ | $S_3$ | $S_4$ | $S_5$ | $S_6$ | $S_7$ | $S_8$ |
|---|---|---|---|---|---|---|---|---|
| $Probe_1$ | a/a | a/a | a/a | c/t | t/t | t/t | t/t | t/t |
| $Probe_2$ | a/a | a/c | a/a | t/t | a/c | a/c | c/t | t/t |

## TABLE 9
### Table for the Class A Occurrences Table

| $Probe_1$ | a/a, 3 | c/t, 1 | |
|---|---|---|---|
| $Probe_2$ | a/a, 2 | a/c, 1 | t/t,1 |

## TABLE 6
### $Probe_1$ A/A and T/T SNPs Occurrences for Fisher Test

| | Class A | Class B |
|---|---|---|
| a/a | 3 | 0 |
| t/t | 0 | 4 |

P-value = 0,0286

## TABLE 7
### $Probe_2$ A/A and C/T SNPs Occurrences for Fisher Test

| | Class A | Class B |
|---|---|---|
| a/a | 2 | 0 |
| c/t | 0 | 1 |

P-value = 0,3333

The formula below gives the exact hypergeometric probability of observing this particular arrangement of data, assuming the given marginal totals, on the **null hypothesis** that Class A and Class B subjects are equally likely to have those SNPs

| Contingency Table | Class A | Class B | |
|---|---|---|---|
| SNP 1 | a | b | a+b |
| SNP 2 | c | d | c+d |
| | a+c | b+d | a+b+c+d (=n) |

Fisher : the probability of obtaining any such set of values is given by the hypergeometric distribution

$$p = \frac{\binom{a+b}{a}\binom{c+d}{c}}{\binom{n}{a+c}} = \frac{(a+b)!\,(c+d)!\,(a+c)!\,(b+d)!}{a!\,b!\,c!\,d!\,n!}$$

# Analysis of results: annotations, dbSNP and PharmGKB links

- Each probe presented in the result window:
  - has a link that allows users to visualize annotations
  - is enriched of links to external databases such as **dbSNP** and **PharmGKB** allowing to automatically retrieve further SNP information
- PharmGKB stores knowledge about the impact of genetic variation on drug response for clinicians and researchers
- dbSNP is a public-domain archive for simple genetic polymorphisms



Data Flow in dbSNP

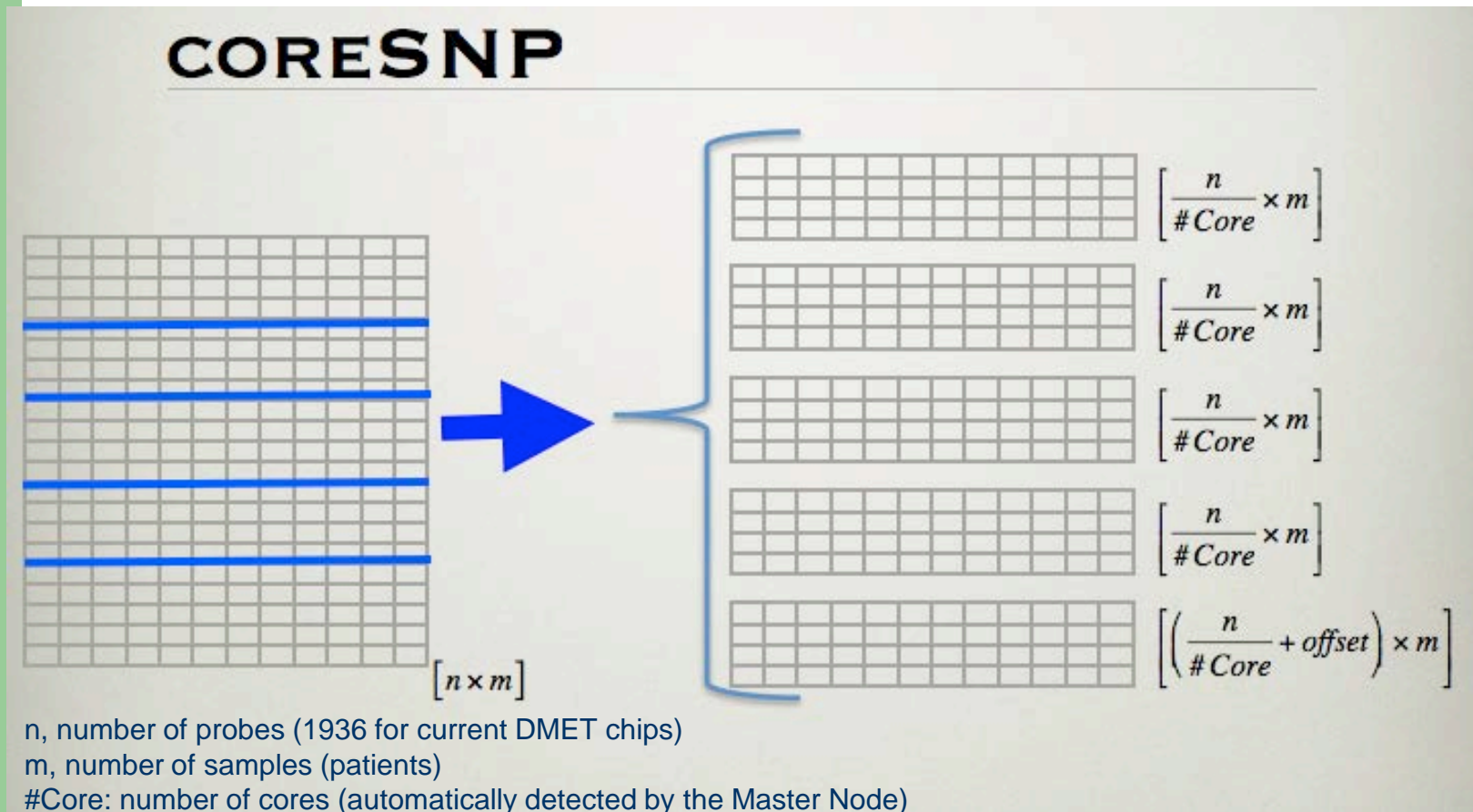| .:: Annotations ::. | |
|---|---|
| Probe Set ID: | AM_10109 |
| dbSNP RS ID: | rs72558192 |
| Chromosome | 10 |
| Physical Position | 96731936 |
| Strand | --- |
| ChrX pseudo-autosomal region 1 | 0 |
| Cytoband | q23.33 |
| Flank | --- |
| Allele A | --- |
| Allele B | --- |
| Associated Gene | ENST00000260682 , exon , 0 , Hs.282624 , CYP2C9 , 1559 , cytochrome P450, family 2, subfamily C, polypeptide 9 ; NM_000771 , CDS , 0 , Hs.282624 , CYP2C9 , 1559 , cytochrome P450, family 2, subfamily C, polypeptide 9 |
| Genetic Map | --- |
| Microsatellite | D10S2360 , downstream , 29989 ; D10S2358 , upstream , 16816 |
| Fragment Enzyme Type Length Start Stop | --- |
| Allele Frequencies | --- |
| Heterozygous Allele Frequencies | --- |
| Number of individuals/Number of chromosomes | --- |
| In Hapmap | --- |
| Strand Versus dbSNP | --- |
| Copy Number Variation | --- |
| Probe Count | --- |
| ChrX pseudo-autosomal region 2 | 0 |
| In Final List | --- |
| Minor Allele | --- |
| Minor Allele Frequency | --- |
| % GC | --- |
| OMIM | |

- A cohort of 19 patients affected by multiple myeloma (MM) treated with aminobisphosphonate zoledronic acid (ZA) were enrolled in a case-control study.
  - 9 patients presented osteonecrosis (ONJ) after the treatment and
  - 10 patients were the control
- The aim of the study was to investigate the association among specific SNPs and the adverse event ONJ induced by ZA.
- Results demonstrated the presence of 8 SNPs that were related to ONJ.
  - Genotypes were determined using DMET Plus GeneChip.
  - Pharmacogenomic profiles were generated by Affymetrix DMET Console software
  - Statistical analysis was performed by two-tailed Fisher's exact test.
- **DMET-Analyzer identified the 8 SNPs** (with same p-values) that were statistically associated with ONJ occurrence, as those identified in the BJH work using manual analysis, in very less time

- To face the increasing volume of genotyping data in pharmacogenomics studies, we designed **coreSNP**, a parallel multi-threaded version of DMET-Analyzer

  – Current DMET chips investigate 1936 probes, each one representing a portion of the genome having a role in drug metabolism (225 ADME-related genes are investigated), but **novel chips may investigate millions probes**

  – **Genome Wide Association Studies (GWAS) involve very large populations of patients, thus datasets for pharmacogenomics studies are increasingly huge**

- **coreSNP** takes into account the data parallelism that can be exploited when analyzing DMET data and uses a simple Master/Slave parallel programming approach to decompose computation

- coreSNP uses a **Fisher Significance (Fs)** threshold to discard probes where SNPs distributions among the two classes A and B are very close, i.e. the Fisher tests involving the SNPs detected on that probe are not computed, thus reducing the computational load of the system:

  – A probe $i$ is discarded if for each SNP $j$, $| FDT[i,j] | <= Fs$, where FDT is a Frequency Difference Table that contains the difference among the frequencies of the SNP j detected on the probe i, respectively in class B and in class A

  – If $Fs=0$, i.e. the most conservative option, only probes with identical alleles distributions in class A and B are discarded.

  – User may also choose to avoid filtering and all probes are tested with Fisher.

# Data Parallelism in coreSNP



n, number of probes (1936 for current DMET chips)
m, number of samples (patients)
#Core: number of cores (automatically detected by the Master Node)

# coreSNP Architecture

# coreSNP: Master and Slave

---

**Algorithm 1** Master Algorithm

---

**Require:** $F_a$, $F_b$, // class A and B data files names
**Require:** $ns_a$, $ns_b$, // number of samples in class A and B
**Require:** $ns$, // total number of samples $ns = ns_a + ns_b$
**Require:** $np$, // number of probes, for DMET data $np = 1936$
**Require:** $ncore$, // number of cores
**Require:** $Fs$, // Fisher Significance Threshold, default 0
**Require:** $Ft$, // Fisher Filter Threshold, default 0.05

1: Begin
2: ComputePartitionsLimits(np,ncore, start, end) // $start_i$, $end_i$, $i = 1, ..., ncore$, point to first and last probe (row) of data files partitions

3: **For Each Core** $i = 1, ..., ncore$
4: **COBEGIN**
5: $Result_i$=slave($F_a$, $F_b$, $ns_a$, $ns_b$, $start_i$, $end_i$, $Fs$, $Ft$)
6: **COEND**
7: Result=merge($Result_1$,...,$Result_{ncore}$, $Ft$) // results are merged and ordered by p-value (not significant results have previously been discarded by Slave threads) and SNPs are annotated with URLs to dbSNP and PharmaGKB
8: print(Result) // results are showed to the user ordered by p-value
9: **end.**

---

**Algorithm 2** Slave Algorithm

---

**Require:** $F_a$, $F_b$, // class A and B data files names
**Require:** $ns_a$, $ns_b$, // number of samples in class A and B
**Require:** $start$, $end$, // limits of data strips assigned to Slave
**Require:** $ns$, // total number of samples $ns = ns_a + ns_b$
**Require:** $Fs$, // Fisher Significance Threshold, default 0
**Require:** $Ft$, // Filter Threshold, default 0.05

1: Begin
2: **For Each Probe** $i = start, start + 1, \ldots, (end - start) + 1$
3: Begin
4: $Oa_i$=ComputeAllelesOccurrences($F_a$,$ns_a$,$i$);
5: $Ob_i$=ComputeAllelesOccurrences($F_b$,$ns_b$,$i$);
6: $FDT_i = \frac{Ob_i}{nsb} - \frac{Oa_i}{nsa}$; // Frequency Difference Table
7: DiscardProbes($FDT_i$, $Fs$); // eventually discard probe $i$
8: $Result_i$=ComputeFisherTests($Oa_i$, $Ob_i$, $Ft$); // computes Fisher tests and discards not significant results
9: End;
10: Return($Result$);
11: **end.**

cannataro@unicz.it

CGW 2017, Krakow, 24 October 2017

**U M G**

*Dubium sapientiae initium*

## Datasets characteristics; Table Reports the Number of Probes and Samples, and the Datasets Dimension on Disk

| Dataset | #probes | #samples | size[Kbytes] |
|---------|---------|----------|--------------|
| DMET | 1,931 | 28 | 245 |
| SNP | 1,062,599 | 28 | 136,673 |

### coreSNP Response Times Without Probe Filtering ($Fs = 0.0$)

| Dataset | 1 | 2 | 4 | 8 | 16 |
|---------|---|---|---|---|----|
| DMET | 20,647 | 10,630 | 5,926 | 3,278 | 2,333 |
| SNP | 11,943,086 | 5,813,502 | 2,890,561 | 1,511,745 | 866,140 |

### coreSNP Response Times with a Light Probe Filtering ($Fs = 0.001$)

| Dataset | 1 | 2 | 4 | 8 | 16 |
|---------|---|---|---|---|----|
| DMET | 20,530 | 10,545 | 5,876 | 3,170 | 2,557 |
| SNP | 10,074,622 | 4,926,823 | 2,432,625 | 1,233,423 | 732,328 |

### Input Probes (#IP), Discarded Probes (#DP), Tested Probes (#TP), Executed Fisher Tests (#FT), Significant Fisher Tests (#SFT), when Varying the Fisher Significance Threshold ($Fs$)

| Fs | #IP | #DP | #TP | #FT | #SFT |
|------|-------|-------|-------|-------|------|
| 0.25 | 1,931 | 1,918 | 13 | 89 | 23 |
| 0.20 | 1,931 | 1,904 | 27 | 212 | 37 |
| 0.10 | 1,931 | 1,582 | 349 | 2,704 | 80 |
| 0.05 | 1,931 | 1,367 | 564 | 4,148 | 80 |
| 0.01 | 1,931 | 876 | 1,055 | 5,932 | 80 |
| 0.001 | 1,931 | 876 | 1,055 | 5,932 | 80 |
| 0.0 | 1,931 | 0 | 1,931 | 7,041 | 80 |



Theoretical Speed-up
DMET dataset Speed-up (Fs = 0.001)
DMET dataset Speed-up (Fs = 0.0)
SNP dataset Speed-up (Fs = 0.001)
SNP dataset Speed-up (Fs = 0.0)

Speed-up / Number of Cores

cannataro@unicz.it

*Dubium sapientiae initium*

- Microarray-based Genomics Data
- Automatic preprocessing of gene expression data
  - **micro-CS**: preprocessing and annotation of gene expression data
- Statistical analysis of genotyping data
  - **DMET-Analyzer**: preprocessing and analysis of DMET (Drug Metabolizing Enzymes and Transporters) SNP (Single Nucleotide Polymorphism) data
  - **coreSNP** (multicore implementation)
- Data Mining analysis of genotyping data
  - **DMET-Miner**: Association Rule Mining to extract from DMET data Association Rules able to correlate the contemporary presence of SNPs with patient's conditions (e.g. TOX vs NOTOX)

# Association Rules

- **Association Rules** are used to find frequent patterns, associations, correlations, or causal structures among sets of items in transaction databases.

- **Let X be an item-set, $X \Rightarrow Y$ an association rule and T a set of transactions of a given database**

- **Support** is an indication of how frequently the item-set appears in the database

  - The support value of X with respect T is defined as the proportion of transactions in the database which contains the item-set X, and is indicated as **supp(X)**

- **Confidence** is an indication of how often the rule has been found to be true

  - The confidence value of a rule $X \Rightarrow Y$ with respect to a set of transactions T, is the proportion of the transactions that contains X which also contains Y.

  - Confidence may be interpreted as an estimate of the conditional probability that a transaction having {X} also contains {Y}

$$\mathrm{conf}(X \Rightarrow Y) = \mathrm{supp}(X \cup Y)/\mathrm{supp}(X).$$

- **Goal: Find all rules $X \Rightarrow Y$ with minimum confidence and support**

- Association rules are usually required to satisfy at the same time:
  - a user-specified minimum support and
  - a user-specified minimum confidence.
- Association rule generation is usually split up into two separate steps:
  1. A minimum support threshold is applied to find all frequent item-sets in a database.
  2. A minimum confidence constraint is applied to these frequent item-sets in order to form rules.
- Main algorithms are Apriori and Frequent Pattern-Growth (FP-Growth).

# DMET-Miner

- DMET-Miner transforms a case-control DMET dataset in a transaction database and it is able to mine Association Rules from such dataset

- DMET-Miner is based on a modified version of the Frequent Pattern-Growth (FP-Growth) algorithm

- The modified version of FP-Growth needs to scan only twice the database to build an FP-Tree, a structure based on extended prefix-tree, making it possible to store in a compressed way crucial information about the frequent patterns.

# DMET-Miner

a) User chooses the data file to be loaded

b) User selects the classes for each sample by clicking on data table.

c) Input table is automatically filtered and inverted by DMET-Miner

*Rules are calculated and visualized*

d) Results are shown to the user. He may visualize annotation and links to external databases.

# DMET-Miner Algorithm and Mined Rules

**Require:** A table of GOA annotation as input dataset $D$
**Ensure:** A set of Weighted Association Rules

1: **Data Structure initialization:** $TDB, FPTree, \beta Tree$

2: $TDB \leftarrow convertInput$

3: **for all** $x \in TDB$ **do**
4:    $wminSupp \leftarrow compute(TDB)$
5:    $frequentItemsList \leftarrow compute(wS, wminSupp)$
6:    **if** $wS(x) \leq wminSupp$ **then**
7:      $TDB.remove(x)$
8:      $frequentItemsList.supportUpdate \leftarrow row$
9:    **else**
10:      $frequentItemsList \leftarrow Update(wS(x))$
11:    **end if**
12: **end for**

13: $descendingSorting(frequentItemsList)$
14: $descendingSorting(TDB)$
15: $FPTree \leftarrow frequentItemsList$

16: **for all** $t \in TDB$ **do**
17:    map each item of t on FPTREE
18:    **if** $perfectmatch$ **then**
19:      $Update(wS(x))$
20:    **else**
21:      $nodeCreation(x)$
22:    **end if**
23: **end for**
24: buildFP-Tree()

25: **for all** $(node \in FPTREE)$ **do**
26:    create $\beta Tree$
27:    **repeat**
28:      **if** $\beta node_{freq} < wminSupp$ **then**
29:        $remove(cpbnode)$
30:      **end if**
31:    **until** $\beta Tree = \emptyset$
32:    $mineRules()$
33:    $saveRules()$
34: **end for**

**Rules of Class NONRESP**

IF( AM_11872_C/C:36 && AM_14348_C/C:36 && AM_13342_G/G:34 && AM_14106_A/A:33 && AM
AM_10188_G/G:25 && AM_10191_C/C:25 && AM_10193_G/G:25 && AM_13141_G/G:24 && AM_103
Sensitivity: 35.0, Specificity: 46.42857142857143
IF( AM_11872_C/C:36 && AM_14348_C/C:36 && AM_13342_G/G:34 && AM_14106_A/A:33 && AM
Sensitivity: 47.368421052631575, Specificity: 70.96774193548387
IF( AM_11872_C/C:36 && AM_14348_C/C:36 && AM_13342_G/G:34 && AM_14106_A/A:33 && AM
&& AM_10191_C/C:25 && AM_10193_G/G:25 && AM_13141_G/G:24 && AM_10330_G/G:21 && AM_
Sensitivity: 29.508196721311474, Specificity: 33.33333333333333
IF( AM_11872_C/C:36 && AM_14348_C/C:36 && AM_13342_G/G:34 && AM_14106_A/A:33 && AM
Sensitivity: 50.0, Specificity: 78.57142857142857
IF( AM_11872_C/C:36 && AM_14348_C/C:36 && AM_13342_G/G:34 && AM_14106_A/A:33
&& AM_12114_G/G:30 && AM_10056_A/A:27 && AM_10177_T/T:25 ) THEN NONRESP conf: 0.92
Sensitivity: 41.66666666666667, Specificity: 60.71428571428571
IF( AM_11872_C/C:36 && AM_14348_C/C:36 && AM_13342_G/G:34 ) THEN NONRESP conf: 0.9
Sensitivity: 58.620689655172406, Specificity: 93.33333333333333
IF( AM_11872_C/C:36 && AM_14348_C/C:36 && AM_13342_G/G:34 && AM_14106_A/A:33 && AM
AM_10188_G/G:25 && AM_10191_C/C:25 && AM_10193_G/G:25 && AM_13141_G/G:24 ) THEN NO
Sensitivity: 38.095238095238095, Specificity: 52.0
IF( AM_11872_C/C:36 && AM_14348_C/C:36 && AM_13342_G/G:34 && AM_14106_A/A:33 ) THE
Sensitivity: 55.00000000000001, Specificity: 89.28571428571429
IF( AM_11872_C/C:36 ) THEN NONRESP conf: 1.0
Sensitivity: 80.0, Specificity: 100.0
IF( AM_11872_C/C:36 && AM_14348_C/C:36 ) THEN NONRESP conf: 1.0
Sensitivity: 67.9245283018868, Specificity: 100.0
IF( AM_11872_C/C:36 && AM_14348_C/C:36 && AM_13342_G/G:34 && AM_14106_A/A:33 && AM
AM_10056_A/A:27 && AM_10177_T/T:25 && AM_10188_G/G:25 ) THEN NONRESP conf: 1.0
Sensitivity: 41.66666666666667, Specificity: 60.71428571428571
IF( AM_11872_C/C:36 && AM_14348_C/C:36 && AM_13342_G/G:34 && AM_14106_A/A:33 && AM
AM_10177_T/T:25 && AM_10188_G/G:25 && AM_10191_C/C:25 ) THEN NONRESP conf: 1.0
Sensitivity: 41.66666666666667, Specificity: 60.71428571428571
IF( AM_11872_C/C:36 && AM_14348_C/C:36 && AM_13342_G/G:34 && AM_14106_A/A:33 && AM
AM_10177_T/T:25 && AM_10188_G/G:25 && AM_10191_C/C:25 && AM_10193_G/G:25 ) THEN NO
Sensitivity: 41.66666666666667, Specificity: 60.71428571428571
Number of Strong Rules founded:13 with confidence value: 60.0 and minimum Support

**Rules of Class RESP**

IF( AM_10330_G/G:52 && AM_12991_C/C:51 && AM_14794_A/A:48 && AM_14525_G/G:45 && AM
Sensitivity: 74.46808510638297, Specificity: 58.53658536585654
IF( AM_10330_G/G:52 && AM_12991_C/C:51 && AM_14794_A/A:48 && AM_14525_G/G:45 && AM
Sensitivity: 75.47169811320755, Specificity: 65.71428571428571
IF( AM_10330_G/G:52 && AM_12991_C/C:51 && AM_14794_A/A:48 && AM_14525_G/G:45 ) THE
Sensitivity: 72.1311475409836, Specificity: 70.37037037037037
IF( AM_10330_G/G:52 && AM_12991_C/C:51 && AM_14794_A/A:48 ) THEN RESP conf: 0.9411
Sensitivity: 72.3076923076923, Specificity: 78.26086956521739
IF( AM_10330_G/G:52 && AM_12991_C/C:51 ) THEN RESP conf: 0.9807692307692307
Sensitivity: 67.56756756756756, Specificity: 85.71428571428571
IF( AM_10330_G/G:52 ) THEN RESP conf: 0.9811320754716981
Sensitivity: 64.55696202531645, Specificity: 88.88888888888889
Number of Strong Rules founded:7 with confidence value: 60.0 and minimum Support v

- DMET-Miner produces rules able to discriminate distinctive features for each class.
- DMET-Miner makes easy to profile in what class new subjects belong (e.g. RESP, NORESP) on the basis of their SNPs

# Comparison among DMET-Miner, Weka and Rapid Miner

- We have compared the FP-Growth algorithm implemented in DMET-Miner with the FP-Growth algorithm implemented in Weka (version 3.6.10) and the FP-Growth algorithm available in RapidMiner (version 5.3.013).

- Weka and RapidMiner were not able to directly load the DMET dataset produced by the DMET platform

- In order to compare DMET-Miner FP-Growth with Weka FP-Growth and RapidMiner FP-Growth on the same conditions, we have given as input to Weka and RapidMiner the filtered dataset produced by our software.

- In this way we ensure that the inputs given to DMET-Miner, RapidMiner and Weka have the same dimensions.

# Datasets

- The comparison of the FP-Growth algorithm implemented in the three software tools was tested using five synthetic DMET datasets. We built the synthetic datasets containing the same number of probes as a real DMET dataset (1, 936 probes) and doubling the number of samples (in the 13 experiments we analyzed five datasets with respectively 25, 50, 100, 200 and 400 samples for each dataset) grouped into two classes.

- We populated these datasets with randomly significantly different distributions of SNPs.

- The datasets contain data related to samples from subjects belonging to two classes: subjects which respond to drugs (class RESP) and subjects which do not respond to drugs (class NONRESP), simulating a classical case-control study.

- The dimensions of the datasets analyzed range from 300KB for the dataset with 25 samples to about 3.1MB for the one with 400 samples

# Fisher Test Filtering

- Number of meaningful probes (rows) after using the Fisher Test Filtering

| Dataset | #Probes | #Probes using Filtering |
|---------|---------|------------------------|
| 25 | 1936 | 45 |
| 50 | 1936 | 48 |
| 100 | 1936 | 57 |
| 200 | 1936 | 76 |
| 400 | 1936 | 78 |

*Dubium sapientiae initium*



Execution time of the algorithms varying the minimum support using the 100 samples dataset. The execution time is obtained in function of the value of minimum support used and the number of mined rules. The dotted line (top part of figure) represents the running time of RapidMiner, the dashed line (middle part of figure), represents running time of Weka, the continuous line (bottom part of figure) represents the running time of DMET-Miner.

Memory consumption of RapidMiner, Weka and DMET-Miner when doubling the size of the input dataset. The dotted line (top part of figure) represents the memory consumption of RapidMiner, the dashed line (middle part of figure), represents the memory consumption time of Weka, the continuous line (bottom part of figure) represents the memory consumption of DMET-Miner.

# OUTLINE

- Experiences at University of Catanzaro in high performance management, preprocessing and analysis of omics data
  - PART I: Genomics data
  - **PART II: Proteomics data**
  - PART III: Interactomics data
- Conclusions

cannataro@unicz.it

*Dubium sapientiae initium*

# Proteomics data

- Introduction to Mass Spectrometry
- MS-Analyzer: service oriented platform for preprocessing and mining of MS data
  - MS-Analyzer = Ontology + Spectra Services + Workflow
- EIPEPTIDI: enhanced protein identification from ICAT MS/MS data

# Early diagnosis of cancer

MECHANISMS OF DISEASE

Mechanisms of disease

## Use of proteomic patterns in serum to identify ovarian cancer

Emanuel F Petricoin III, Ali M Ardekani, Ben A Hitt, Peter J Levine, Vincent A Fusaro, Seth M Steinberg, Gordon B Mills, Charles Simone, David A Fishman, Elise C Kohn, Lance A Liotta

Recently MS has been used to analyse low weight proteins present in serum/plasma for early diagnosis of tumours

# MS-based Proteomics



**Mass Spectrometer**

Sample introduction

Interface to Vacuum | Ion Source | Mass Analyser | Detector

High Vacuum

**Electron Impact Ionisation**

**Raw Mass Spectrum**

**Cleaned Mass Spectrum**

Intensity

Mass to Charge ratio [m/Z]

*Data Preparation*

**Workflow Execution
(i.e. Supervised Classification)**

**Data Pre-Processing**

Subtract Base Line
Normalization
Binning
Peak Identification
Peak Alignment

**Knowledge Models**

**Preprocessed Cleaned Spectra Database**

cannataro@unicz.it

CGW 2017, Krakow, 24 October 2017

# MS at UMG

MALDI-TOF



Applied Biosystems Voyager DE-STR

ESI-QqTOF



Applied Biosystems QSTAR XL LC-MS/MS

# Tandem mass spectrometry (MS/MS)

MS1: MS spectrum
from 0 to 1000 m/z (Da)

MS2: MS/MS
on m/z 850.1

Peak selection



cannataro@unicz.it

# Spectra Preprocessing

- *Noise Reduction*
  - *Base line subtraction* flattens the base profile of a spectrum
  - *Smoothing* reduces the noise level in the whole spectrum.
- *Binning*
  - groups measured data into bins
    - aggregate intensity (e.g. the sum of the intensities in the bin)
    - representative m/Z value (e.g. the median or the one with maximum intensity)
- *Normalization*
  - aims to make intensity comparable across different spectra.
- *Peaks alignment*
  - The same peak (e.g. the same peptide) may have different m/Z values across samples
  - finds a common set of peak locations (i.e. m/Z values) in a set of spectra, so that all spectra have common m/Z values for the same biological entities

# Spectra analysis

- Biomarker discovery
  - Analysis of collection of spectra (usually by using data mining or machine learning) to find discriminating peaks among different conditions (e.g. healthy, diseased)

- Protein/Peptide identification
  - Identification of protein/peptide associated to a peak, e.g. through MS/MS and database search (Mascot, Sequest, X!Tandem)

- Qualitative vs quantitative proteomics
  - MALDI-TOF or SELDI-TOF give a snapshot of a sample, without providing quantitative information nor complete identification
  - Labelling techniques, such as ICAT or SILAC, coupled to tandem MS (E.g. ICAT-based LC-MS/MS) allow quantification and identification of proteins

# MS data

- The basic MS data is a long sequence of (intensity, m/z) value pairs
- Indeed, different MS data formats do exist depending on:
  - Ionization (ESI, MALDI, TOF)
  - Separation (e.g. GAS or Liquid Chromatography)
    - More spectra at different (retention) times
  - MS approach (e.g. MS vs MS/MS)
    - One spectrum/sample vs more spectra/sample
  - Labeling
    - intensity in a sample is a measure relative to a control one
- Moreover, each MS company usually uses a proprietary format
- Heterogeneity in spectra data is an issue

*Dubium sapientiae initium*

# mzData

- mzData was the first XML-based data model defined by HUPO-PSI (Human Proteome Organization-Proteomics Standard Initiative) to standardize mass spectrometry-based experimental data
- mzData schema comprises a
  - **description** element describing metadata about the experiment (e.g. administrative information, data about instrument and software used to generate the spectra),
  - and a **spectrumList** element that encloses a set of **spectrum** elements.
- Each spectrum element stores, respectively, all the **m/z** and **intensities** values of the spectrum as Base64 encoded strings.
  - Base64 encoding allows to represent arbitrary sequences of octets through a 65-character alphabet (each 6 bits are represented as a printable character).
  - Base64 allows to easily transmit MS data over Internet protocols and can be stored into XML documents.

# Issues in Mining Spectra Data

- **Loading** of the raw spectra produced by mass spectrometer
- **Converting** (eventually) in a standard format (mzXML, mzData),
- **Preprocessing** of the raw spectra data,
- **Preparation** of the data mining input file (e.g. Weka ARFF file),
- **Data Mining analysis** (e.g. classification) of mass spectra,
    - CRISP-DM methodology,
    - Clementine, Weka, and SPSS platforms
- **Knowledge Models visualization** (e.g. decision tree)
    - represented in a standard language as PMML

- **Key issues**
    - Service Oriented Distributed Computation
    - Efficient spectra management
    - Provenance data

*Dubium sapientiae initium*

# MS-Analyzer

- MS-Analyzer is a software platform for the preprocessing, management and data mining analysis of MS data. It uses
  - domain ontologies to model
    - software tools (e.g. preprocessing and mining) and their relationships
    - data sources (e.g. MALDI-TOF, LC-MS/MS spectra datasets)
    - constraints (e.g. binning cannot be applied twice)
  - and workflow techniques to design "in silico" experiments.

- MS-Analyzer:
  - acquires, preprocesses and manages MS data
  - offers filtering, preprocessing, and DM services
  - sharing experiments data, workflows and knowledge

**Ontology-based Workflow Designer**

**WF Editor**
- composition
- browsing
- selection
- visualization

**Ontology Assistant**
- browsing
- querying

Ontologies

**Ontology manager**

Resource Discovery Services

**UDDI/MDS**

Metadata WSDL

WF Schema Abstract, Concrete WF

**Workflow Manager**

WF Translator

WF Browser

WF Scheduler

WF Monitor

**SELDI MALDI ICAT ...**

$WS_1$
$WS_2$
**Spectra Management Services**

$WS_1$
$WS_2$
**Spectra Preprocessing Services**

$WS_1$
$WS_2$
**Spectra Preparation Services**

**Grid Middleware**

$WS_1$
$WS_2$
**Data Mining Services**

$WS_1$
$WS_2$
**Spectra Visualization Services**

**SpecDB APIs**

**RSR**
raw MS spectra

**PSR**
pre-processed MS spectra

**PPSR**
prepared pre-processed MS spectra

cannataro@unicz.it

17, Krakow, 24 October 2017

# MS-Analyzer Ontologies

# MS-Analyzer GUI

**U M G**

*Dubium sapientiae initium*



Design Sheet

Dataset Area

Log Area

Ontology Area

View Area

cannataro@unicz.it

CGW 2017, Krakow, 24 October 2017

# Ontology Manager

- Allows
  - loading user-defined OWL ontologies
  - browsing and searching bioinformatics tools
  - green services can be used in the editor by drag&drop

# Dataset manager

*Dubium sapientiae initium*

- Each experiment is represented as a tree containing
  - raw and preprocessed datasets
- Different source spectra are supported
  - peak list, TXT, CSV, ARFF
- Different formats to be given in input to data mining
  - ARFF, Clementine
- The user can
  - Add a new experiment
  - Loading of experiment

# Workflow Editor

- UML-based notation
- Services are taken from the ontology pane
- Datasets are taken from the Dataset pane
- Constraints expressed by the ontology are enforced at composition time
- Ongoing: generation of a BPEL WF schema to be scheduled by a (Grid) workflow engine



cannataro@unicz.it

# A Grid Environment for High-Throughput Proteomics

Mario Cannataro*, *Associate Member, IEEE*, Annalisa Barla, Roberto Flor, Giuseppe Jurman, Stefano Merler, Silvano Paoli, Giuseppe Tradigo, Pierangelo Veltri, and Cesare Furlanello, *Associate Member, IEEE*

**MS-Analyzer on the Grid**

# Protein Identification with Tandem Mass Spectrometry (MS/MS)

- Detection: first MS detects and selects the most abundant peptides (detection)

- Discovery: second MS analyzes such peptides producing a new spectrum

- Identification/Quantitation: the list of peaks is used to query different protein/peptide identification tools

  - MASCOT, ProteinProspector, Sequest

- To obtain a **quantitative** measurement about the protein abundance, MS may be coupled with opportune sample preparation protocols.

  - e.g. proteins can be identified by using tandem mass spectrometry with ICAT (Isotope Coded Affinity Tag) protocol for sample preparation.

# ICAT + MS/MS protein identification

**Control Sample** → **Label with Light ICAT Reagent**

**Test Sample** → **Label with Heavy ICAT Reagent**

↓

**Combination of Control and Test**

↓

**Digestion Isolation of Cys-Labeled Peptides**

↓

**Mass Spectrometry**

→ Pro ICAT

Swiss Prot DB

*Identified and Quantified Proteins*

Table 2. Single Sample Dataset

| Protein ID | H:L Ratio |
|---|---|
| gi112911 | 0.7554 |
| gi123508 | 0.5633 |
| gi135807 | 1.1732 |
| gi12507 | 0.2620 |
| gi1708182 | 0.6928 |
| gi116117 | 1.1939 |
| gi122801 | 1.0241 |
| gi122910 | 0.9758 |
| gi139653 | 1.0792 |

Table 3. Multiple Samples Dataset

| Protein ID | Sample1 | Sample2 | Sample3 |
|---|---|---|---|
| gi125507 | 0.2620 | 1.9812 | 1.4462 |
| gi1708182 | 0.6926 | 1.0282 | 1.0625 |
| gi116117 | 1.1939 | 1.2932 | 1.8219 |
| gi122801 | 1.0241 | 1.5374 | 0.9964 |
| gi112910 | 1.0792 | 1.3566 | 1.1751 |
| gi139653 | 1.0792 | 1.3566 | 1.1751 |
| gi46577680 | 1.0770 | 1.2592 | 1.1103 |
| gi113936 | 1.0035 | 1.5346 | 0.8608 |

cannataro@unicz.it

CGW 2017, Krakow, 24 October 2017

*Dubium sapientiae initium*

# EiPeptiDi: Enhanced ICAT Peptide Identification Discovery Tool

- Identified peaks (m/Z, retention time, intensity) + name of originating peptide/protein are stored in a database
- In new experiments, **peaks detected by the first MS, but not selected for second MS** can be identified by searching for "similar" peaks (m/Z, retention time, intensity)  in the database

Swiss Prot DB

Pro ICAT - Protein Identification and Quantification Module

EIPeptiDi Tool Boost Discovery Module

MS/MS routine

| | Pept. *a* | Pept. *b* | Pept. *c* | Pept. *d* |
|---|---|---|---|---|
| Sample 1 | 0.75 | 0.72 | | |
| Sample 2 | | 1.03 | | |
| Sample 3 | 0.95 | | 0.82 | |
| Sample 4 | | | | 0.79 |
| Sample N | 1.12 | | | |

| | Pept. *a* | Pept. *b* | Pept. *c* | Pept. *d* |
|---|---|---|---|---|
| Sample 1 | 0.75 | 0.72 | 0.81 | 0.97 |
| Sample 2 | 0.89 | 1.03 | | 0.71 |
| Sample 3 | 0.95 | 1.23 | 0.82 | |
| Sample 4 | 0.88 | | 1.33 | 0.79 |
| Sample N | 1.12 | 0.97 | 0.86 | 1.03 |

CGW 2017, Krakow, 24 October 2017

# EiPeptiDi Algorithm



**Sample S1:**

2 peptides identified and quantified

**Sample S2:**

2 peptides quantified, but **only one identified**

**EIPEPTIDI** cross validates peptides identification across the data set

**Retention time windows**

$$St_i + delta <= rt_x <= St_i - delta$$

**Mass windows**

$$(Sm_i - (m_x * deltaMz))$$
$$<= MwF <=$$
$$(Sm_i + (mx * deltaMz))$$

cannataro@unic

**procedure** $PeptideDiscovery(F, NF)$
// F contains the peptides found within samples with masses, retention times
// NF contains masses and retention times not assigned within samples
**const** $MAX\_MT = 0.00003$; //mass tolerance 30 ppm
**const** $MAX\_RTT = 3$; //retention time tolerance
**var** $\Delta_m, \Delta_{StartTime}, \Delta_{EndTime}$: **real**;
**begin**
    **for each** $f_i = (s_i, p_i, r_i, m_i) \in F$ **do begin**
        **for each** $nf_j = (s_j, r_j, m_j) \in NF$ **do begin**
          //calculate $m_j$ tolerance
          $\Delta_m := MAX\_MT * m_j$;
          $\Delta_{StartTime} := \mathbf{abs}(r_j.startTime - r_i.startTime)$;
          $\Delta_{EndTime} := \mathbf{abs}(r_j.endTime - r_i.endTime)$;
          // Verify Mass And Retention Time
          **if** $((m_j - \Delta_m < m_i < m_j + \Delta_m)$ **and**
             $\Delta_{StartTime} <= (MAX\_RTT/2)$ **and**
             $\Delta_{EndTime} <= (MAX\_RTT/2))$ **then**
           // Assign the peptide $p_i$ to not found $nf_j$
           $nf_j.AddPeptideAsIdentified(p_i)$;
           //$p_i$ is also contained in sample $s_j$
        **end**;
      **end**;
**end** PeptideDiscovery;

# EiPeptiDi GUI



EiPeptiDi is available on line at:

http://bioingegneria.unicz.it/~veltri/projects/eipeptidi

cannataro@unicz.it

*Dubium sapientiae initium*

- In a 7 samples data sets we increase the number of Identified Peptides of 50 % (av) using EiPeptiDi.



Data overlap increases

# Experimentation

- EIPETIDI has been successfully used to improve the identification of proteins used to discriminate between different classes of familial adenomatous polyposis (FAP) patients

    – Barbara Quaresima, Telma Crugliano, Marco Gaspari, Maria Concetta Faniello, Paola Cosimo, Rosa Valanzano, Maurizio Genuardi, Mario Cannataro, Pierangelo Veltri, Francesco Baudi, Patrizia Doldo, Giovanni Cuda, Salvatore Venuta, Francesco Costanzo, A proteomics approach to identify changes in protein profiles in serum of familial adenomatous polyposis patients, Cancer Letters 272(1):40-52,8 December 2008, http://dx.doi.org/10.1016/j.canlet.2008.06.021

# OUTLINE

*Dubium sapientiae initium*

- Experiences at University of Catanzaro in high performance management, preprocessing and analysis of omics data

  - PART I: Genomics data

  - PART II: Proteomics data

  - **PART III: Interactomics data**

- Conclusions

# Interactomics Data

- Interactomics
  - Protein-to-Protein Interactions (PPI)
  - PPI Databases and Standards
  - Protein Complex Prediction

- Experiences at University of Catanzaro:
  - **IMPRECO**: a meta-predictor for protein complexes
  - **CytoMCL**: Markov Clustering of Metabolic Networks in Cytoscape
  - **CytoSeVis**: Semantic Similarity-Based Visualization of PPI Networks in Cytoscape
  - **ONTOPIN**: using ontologies for querying PPI databases

# Interactomics

- **Interactomics** is the study of the whole set of protein interactions
  - Protein to Protein Interaction (PPI) network.
- PPI Networks are usually modeled as undirect graphs.
  - Global analysis investigated main properties.
- Different network models have been proposed:
  - Scale Free
  - Random Graphs
  - Geometric Random Graph

Dubium sapientiae initium

# Flow of Information in Interactomics



From wet lab to in silico analysis and prediction

*Dubium sapientiae initium*

- Verified Interactions DBs store experimentally determined (Yeast 2 Hybrid, Protein Micro Arrays, Mass Spectrometry) PPI data
  - **DIP** (Database of Interacting Proteins), **BIND** (Biomolecular Interaction Network Database), **MIPS**, **MPCDB, CORUM**
- Predicted Interactions DBs store interactions predicted in silico:
  - **OPHID** maps experimental interactions determined in model organisms into human interactions.
  - **POINT** projects verified interactions into human orthologs and filters interactions considering functional information.
  - **IntNetDB** predicts interaction by integrating different information (mRNA, co-expression, sequence similarity)
- Standards
  - BioPax, IMEx, PSI-MI, SBML, SIF

*Dubium sapientiae initium*

# Complex Prediction Algorithms

- A **protein complex** is a group of two or more associated proteins which interact sharing the same biological goal
  - i.e. a cluster in the PPI graph

- The Markov Cluster algorithm (MCL)  [van Dongen 2000] finds clusters on a graph by simulating a stochastic flow and then analyzing its distribution
- The MCODE algorithm, takes in input an interaction network and tries to find complexes by building clusters.
- The  Restricted Neighborhood Search (RNSC) predictor, after an initial random clustering, uses  a cost-based local search algorithm based on the tabu heuristic



cannataro@unicz.it

# MCL



Input Network

Clustered Network

# IMPRECO: A Tool for Improving the Prediction of Protein Complexes

**U M G**
*Dubium sapientiae initium*

**Mario Cannataro, Pietro H. Guzzi, Pierangelo Veltri**

**Bioinformatics Laboratory, University "Magna Græcia" of Catanzaro, Italy**
**cannataro@unicz.it**

# IMPRECO: Improving the prediction of protein complexes

# The IMPRECO Algorithm

- The rationale underlying the proposed meta-predictor is to combine different predictor results using an integration algorithm able to gather (partial) results from different predictors

- The integration algorithm starts by integrating results (i.e. clusters) obtained by running different available predictors.

- Three different cases are considered by evaluating the topological relations among clusters coming from the considered predictors:

  1. **equality:** the same clusters are returned by all (or by a significant number of) predictors,

  2. **containment:** it is possible to identify a containment relation among (a set of) clusters returned by all (or by a significant number of) predictors;

  3. **overlap:** it is possible to identify an overlap relation among (a set of) clusters returned by all (or by a significant number of) predictors;

cannataro@unicz.it

*Dubium sapientiae initium*

# Validation of Predicted Complexes



- To estimate the integration quality, IMPRECO uses an **evaluation module based on a reference database, i.e. a catalog of verified complexes** (e.g. the MIPS catalog)

- For each cluster the evaluation module calculates the measurements of sensitivity, positive predictive values and accuracy .

  - The first measure is an average representing the fraction of proteins of a complex that are found in a common cluster. When only a big cluster is found, the sensitivity tends to one.

  - The second measure represents the fraction of members of a cluster that belong to a given complex. When each protein belongs to one cluster, PPV is 1, conversely to the previous measure.

  - Thus, the third measure, being the geometric average of sensitivity and ppv represents a trade-off.

# The IMPRECO TOOL

# Preliminary Results/1

In the first experiment we used a yeast network of 1094 nodes and 14658 edges and we run the MCL, RNSC and MCODE algorithms, obtaining respectively 165, 306 and 73 clusters.

| Membership | Strong ▼ | Value | 0 |
|---|---|---|---|
| Sub Graph Intersection | Smallest ▼ | | |
| OverLap Intersection | Smallest ▼ | OverLap Percentage | 0.1 |
| Integration Schema | Templat... ▼ | | |
| Threshold of Dimension | 3 | | |

| Parameter | Value |
|---|---|
| TM | 1 |
| IS | Biggest |
| II | Biggest |
| Overlap Percentage | 0.1 |
| TD | 2 |

| Parameter - Algorithm | MCODE | MCL | RNSC | IMPRECO |
|---|---|---|---|---|
| Clusters No. | 73 | 165 | 306 | 155 |
| Sensitivity | 0,34 | 0,47 | 0,46 | 0,89 |
| PPV | 0,48 | 0,69 | 0,59 | 0,59 |
| Accuracy | 0,40 | 0,57 | 0,52 | 0,70 |

# CytoMCL:
# Markov Clustering of Metabolic Networks in Cytoscape

Mario Cannataro, Pietro H. Guzzi

Bioinformatics Laboratory,

University "Magna Græcia" of Catanzaro, Italy

cannataro@unicz.it

U M G
*Dubium sapientiae initium*

# CytoMCL rationale

- The Markov Clustering Algorithm (MCL) is a well-known algorithm for clustering graphs.
- Cytoscape is a tool for visualizing and analyzing networks based on an extensible architecture.
- Nevertheless MCL does not provide a graphical user interface and cannot be used in the Cytoscape platform
- CytoMCL is a Cytoscape plugin that finds clusters in a graph by using the Markov Clustering Algorithm.
  - Based on an intuitive interface it is able to load a network from Cytoscape, to analyze it and to visualize resulting clusters into Cytoscape.

# Starting the plugin

# CytoSeVis: Semantic Similarity-Based Visualization of PPI Networks

**Mario Cannataro**, Pietro Hiram Guzzi

University "Magna Graecia" of Catanzaro, Italy

cannataro@unicz.it

# PINs Visualization

- PINs visualization tools should:
  - handle high-dimensional PINs
  - meet users' requests in terms of low response time, interactivity, ease of use
  - support some form of graph analysis



- PINs visualization tools comprise:
  - Collections of **layout algorithm** (e.g. circular, tree, hierarchical, Force-directed);
  - Different graphical **rendering algorithms** (2D, 3D)
  - **Graph-analysis** and **annotation** (clustering, statistics)

- Main tools include Cytoscape and NAViGaTOR

Agapito et al. Visualization of protein interaction networks: problems and solutions, BMC Bioinformatics 2013

Pavlopoulos et al. BioData Mining, 2008

# Visualization in Cytoscape

- **Layouts**: Spring-Embedded, Circular, Grid, Force-Directed and Organic;
- **Rendering**: 2D
- **Network analysis**: basic statistical analysis:
  - through the plug-in manager it is possible to add external modules
- **Cytoscape does not support full semantic similarity analysis**
  - GOlorize provides class-guided network visualization (used with BINGO)
  - BINGO finds GO categories overrepresented in a selected part of the network

cannataro@unicz.it

CGW 2017, Krakow, 24 October 2017

*Dubium sapientiae initium*

# Biological Knowledge is encoded into many Biological Ontologies

# Semantic Similarities

- Biological entities and molecules are also associated with GO terms representing their functions, biological roles and localization.

- The association of GO terms and biological molecules is called **annotation process** and it may be performed manually under the supervision of an expert, or automatically.

- There exist currently 17 different annotation processes that are identified by an **evidence code.**

- The **set of annotations** is stored in different databases, such as *the Gene Ontology Annotation Database* (GOA)

Guzzi et al, Briefings in Bioinformatics, 2012

*Dubium sapientiae initium*

**TP11**

## Annotations

| Genes Proteins | | GO TERMS | GO:0008152, GO:0004807, GO:0016853, GO:0003824 |
|---|---|---|---|

| Experimental Evidence Codes | | Computational Analysis Evidence Codes | |
|---|---|---|---|
| EXP | Inferred from Experiment | ISS | Inferred from Sequence or Structural Similarity |
| IDA | Inferred from Direct Assay | ISO | Inferred from Sequence Orthology |
| IPI | Inferred from Physical Interaction | ISA | Inferred from Sequence Alignment |
| IMP | Inferred from Mutant Phenotype | ISM | Inferred from Sequence Model |
| IGI | Inferred from Genetic Interaction | IGC | Inferred from Genomic Context |
| IEP | Inferred from Expression Pattern | RCA | Inferred from Reviewed Computational Analysis |
| **Author Statement Evidence Codes** | | **Curator Statement Evidence Codes** | |
| TAS | Traceable Author Statement | IC | Inferred by Curator |
| NAS | Non-traceable Author Statement | ND | No biological Data available |
| **Automatically-assigned Evidence Codes** | | **Obsolete Evidence Codes** | |
| IEA | Inferred from Electronic Annotation | NR | Not Recorded |

# Ontology-based analysis

- GO and GO annotations enable the use of a set of **analysis methodologies** that have been **defined for the ontologies.**

- There exists, in fact, a set of analysis methodologies of the terms of ontologies as well as of the annotated entities

- An important **tool for ontology-based analysis** is represented by the **pairwise semantic similarities measures (SSM ).**
    - **SSM** are mathematical functions that quantify the similarities of two terms belonging to an ontology into a numerical value.
    - **SSM** are usually **defined for the comparison of two terms** but they can be easily extended to consider group of terms as input.

# Semantic Similarity

- In computational biology, semantic similarity measures are generally based on Gene Ontology and are often used to compare both GO terms and gene products.

- There exist currently many available semantic similarity measures as listed in Pesquita et al. that can be categorized using different parameters on the basis of the steps they employ to determine the semantic value.

Dubium sapientiae initium

*Dubium sapientiae initium*

- **Pairwise: Two Terms**
- **Groupwise: Set of terms, i.e. proteins annotated with a set of terms**

**Proteins are described with a set of annotating terms**

TPI1 → GO:0008152,GO:0004807,GO:0016853, GO:0003824

PGAM2 → GO:0006096,GO:0008152,GO:0016853, GO:0003824,GO:0016868

| MEASURE | METHOD |
|---|---|
| Resnik | Pairwise |
| ResnikGraSM | Pairwise |
| Lin | Pairwise |
| LinGraSM | Pairwise |
| JiangConrath | Pairwise |
| JiangConrathGraSM | Pairwise |
| Relevance | Pairwise |
| Kappa | Groupwise |
| Cosine | Groupwise |
| SimGIC | Groupwise |
| CzekanowskiDice | Groupwise |

# Cyto-SeVis

- CytoSeVis is a Cytoscape plugin that is able to visualize protein interaction networks in a semantic similarity space

- CytoSevis is based on an intuitive interface and is able to load a network from Cytoscape

- It  loads the semantic similarities provided as separate files and visualizes resulting coloured network into Cytoscape workspace

# CytoSeVis Availability

- CytoSeVis for Windows, Linux and Mac OSX platforms is available under GPL license.

- Download at Cytoscape web site: http://cytoscape.org

- Semantic similarities for yeast available at:http://bioingegneria.unicz.it/~guzzi/ss



cannataro@unicz.it

# CytoSevis Operation

Once that a network has been loaded.

1. User has to select the CytoSeVis plugin and an initial node.

2. Then user has to select and load the similarity file

# CytoSeVis SS input file

SS input file is a NxN precomputed matrix containing SS among each couple of proteins

SS are calculated with csbl.go R package: http://csbi.ltdk.helsinki.fi/csbl.go/ (Ovaska et al, BioData Mining, 2008)

**Available SS:**

Cosine

Czekanowski

Dice

JiangConrath

Kappa

Lin

LinGrasm

Relevance

Resnik

Weighted Jaccard

Resnik Grasm



FileStampato3.txt - Blocco note

File   Modifica   Formato   Visualizza   ?

```
NA   node0 node1 node2 node3 node4 node5
node0 NA   0.17 0.04 0.07 0.1  0.14
node1 0.17 NA   0.24 0.27 0.3  0.34
node2 0.04 0.24 NA   0.44 0.47 0.5
node3 0.07 0.27 0.44 NA   0.64 0.67
node4 0.1  0.3  0.47 0.64 NA   0.84
node5 0.14 0.34 0.5  0.67 0.84 NA
```

**UMG**

*Dubium sapientiae initium*

# CytoSeVis Use



- Semantic similarities are translated into different colors.
  - The white node is the reference seed node.
  - All the other similarities are calculated with respect to this node.
- In the example, a darker color represents a lower similarity.

- Semantic similarity adds a new dimension for PINs visualization and analysis
- CytoSeVis is a Cytoscape plugin for semantic similarity-based visualization of PINs
  - http://cytoscape.org
- Semantic similarities are available at:http://bioingegneria.unicz.it/~guzzi/ss
- See also:
  - Guzzi et al Briefings in Bioinformatics 2012,
  - Guzzi et al IEEE BIBMW 2011, ICIAP 2011

# OUTLINE

- Experiences at University of Catanzaro in high performance management, preprocessing and analysis of omics data

  - PART I: Genomics data

  - PART II: Proteomics data

  - PART III: Interactomics data

- **Conclusions**

*Dubium sapientiae initium*

*Dubium sapientiae initium*

- **DMET-Analyzer** supports the automatic statistical analysis in DMET-based pharmacogenomics data and it is able to find statistically relevant subsets of SNPs that separate two input classes

- **DMET-Miner** extracts Association Rules from DMET-based pharmacogenomics data

- **OSAnalyzer** analyzes genomics data annotated with clinical data and computes the OS and PFS curves related to the presence/absence of SNPs in the population under investigation

# Acknowledgements

*Dubium sapientiae initium*

➢ **Prof. Mario Cannataro**

➢ **Prof. Pierangelo Veltri**

➢ **Prof. Pietro H. Guzzi**

➢ Dr. Agapito Giuseppe, PhD, PostDoc;

➢ Dr. Calabrese Barbara, PhD, PostDoc;

➢ Dr. Giuseppe Tradigo, PhD, PostDoc;

➢ Dr. Vizza Patrizia, PhD, PostDoc;

➢ Ing. Cristiano Francesca, PhD Student;

➢ Ing. Milano Marianna, PhD Student;

➢ Ing. Mirarchi Domenico, PhD Student;

➢ Dr. Chiara Zucco, PhD Student;

➢ Ing. Pietro Cinaglia, PhD Student.

cannataro@unicz.it