# Bioinformatics and Computerscience

# Systems Biology

| Data collection | Network Inference | Network-based dataintegration |
|---|---|---|

1. ARRAY BASED

2. NEXT-GEN SEQUENCING
RNA-Seq analysis
ChIP-seq
Bulked segregant analysis

1. Sequence-based data analysis
MotifSuite
ModuleDigger
Crossed

2. Network reconstruction
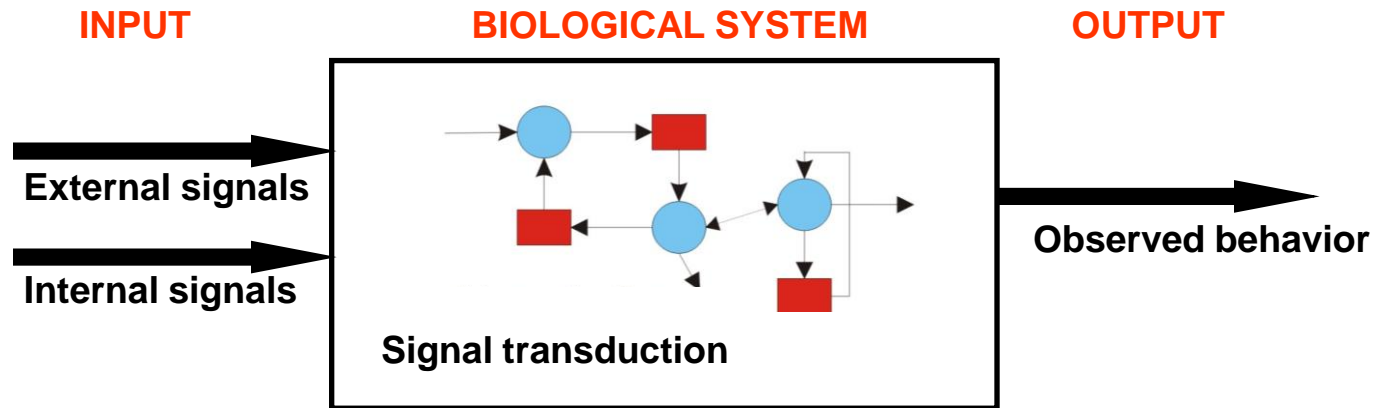Lemone
Distiller
Comodo
Bayesian network reconstruction

1. Network-based analysis of unstructured gene lists
2. Network-based gene prioritization
3. Network-based eQTL analysis
4. Network-based subtyping

- Development of methods that assist systems biologists
- Methods based on data-mining, statistics
- Unique in combining biologically relevant assumptions with rigorous statistical and datamining framework (pragmatic but not too much ad hoc)
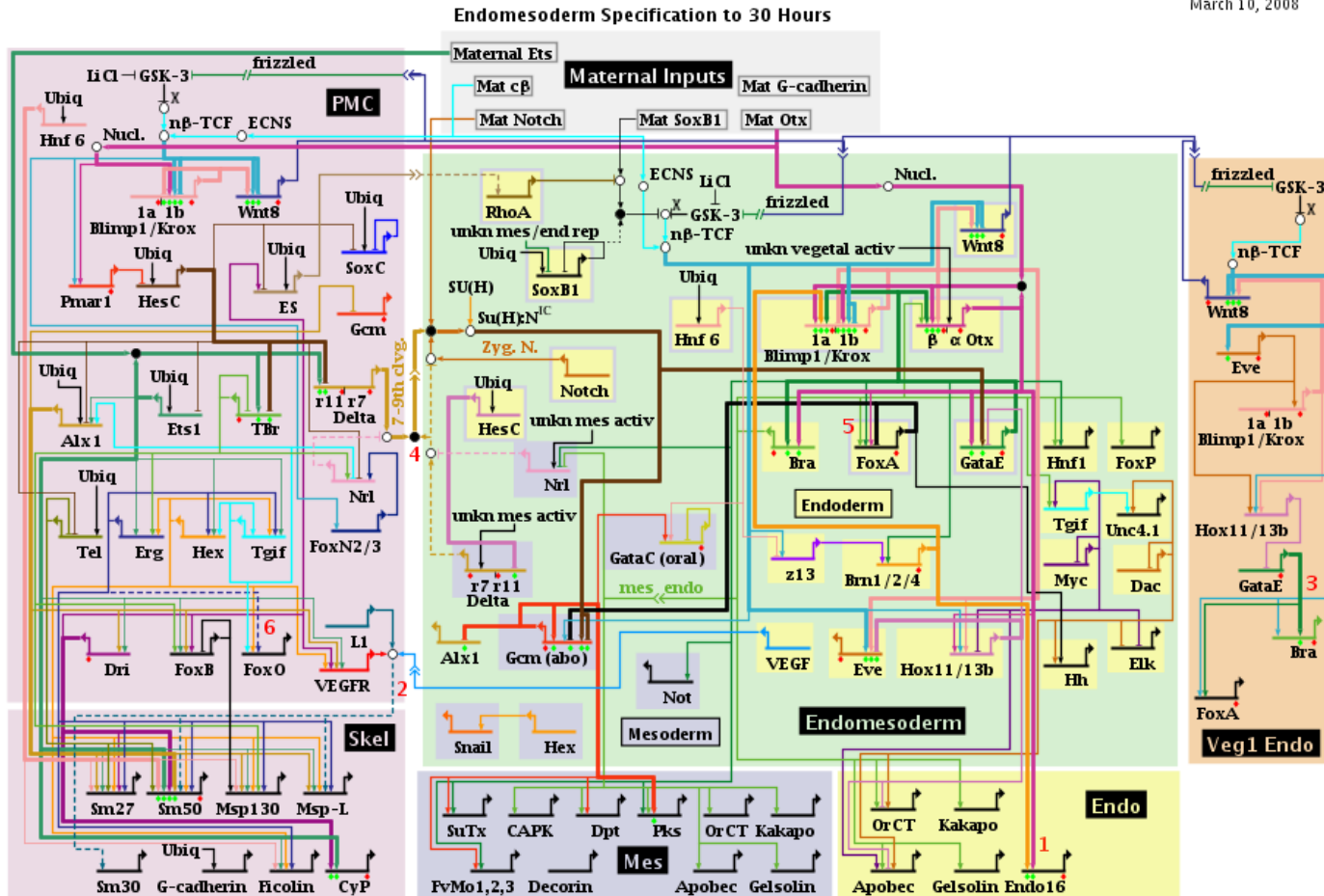- All tools have been validated on real biological cases

**http://bioinformatics.psb.ugent.be/DBN/dbn/software**

UNIVERSITEIT GENT

KU LEUVEN

# Systems Biology



**INPUT**

**BIOLOGICAL SYSTEM**

**OUTPUT**

**External signals**

**Internal signals**

**Signal transduction**

**Observed behavior**

# Systems Biology



Endomesoderm Specification to 30 Hours
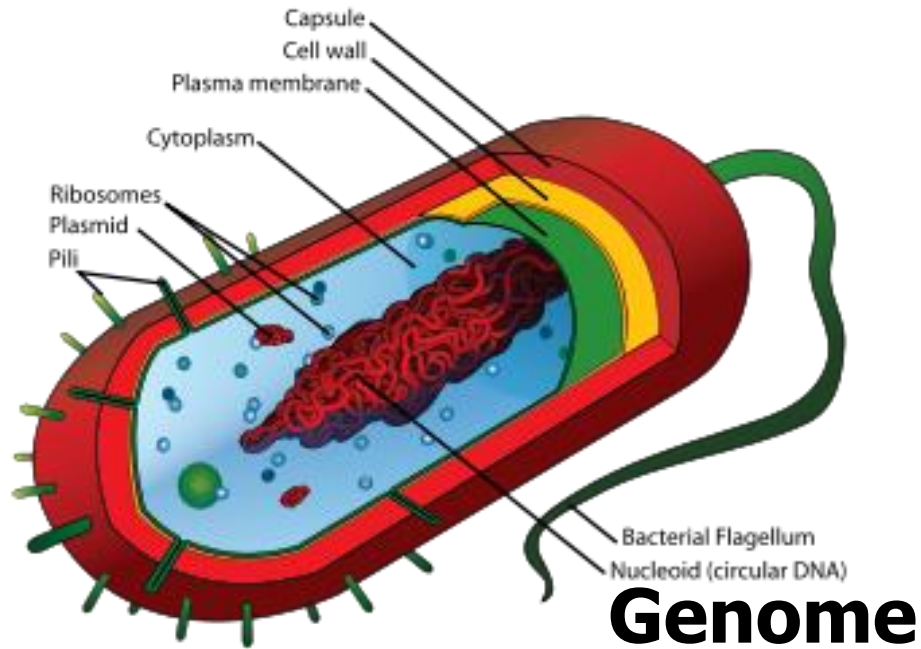
March 10, 2008

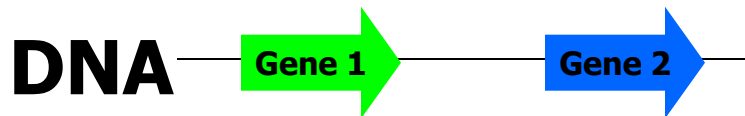Copyright © 2001–2008 Hamid Bolouri and Eric Davidson

Ubiq=ubiquitous; Mat = maternal; activ = activator; rep = repressor;
unkn = unknown; Nucl. = nuclearization; χ = β-catenin source;
nβ-TCF = nuclearized b-β-catenin-Tcf1; ES = early signal;
ECNS = early cytoplasmic nuclearization system; Zyg. N. = zygotic Notch

# Systems Biology



Capsule
Cell wall
Plasma membrane
Cytoplasm
Ribosomes
Plasmid
Pili
Bacterial Flagellum
Nucleoid (circular DNA)

**Genome**

# Systems Biology

**DNA**



AGCACTGTCCACTGCATGGTGAGGATGGGGGTGAGCTCCCT
TTGTGGCTAGGTGCTTAAACGTCTATCGGACGCTCA**GTGAA**
**GGGCTATTGCGGGGTAACAGGAAACCCTGAGGTGCT**
**GATAGGTCAAAGATGGAGAAGGCGTAGGGATATGTG**
**TTGGCAGAGGAACCGAAGAATACCAGGCCATTTCCG**
**AAGCCACTCATTTTCCTCGAATTCCTTCTTTAT**GCCTTC
AGTCTCTATTGACCGTAAATTTGGTTGTTGTCTCCCAGCTGT
TTATTTCTGTAACAGATCTTGGAGGCTGCGGTCTGGATCCCT
CGCCAAGAACCAGATCCAGGAGAAACGTGCTCAACGTGCA
GCTCTGCTCCTACTGATTATAGCCCCACAGATGACATCGCTC
CATAGTCACACCAAGTCTCCTGTGGGAGTCTTGCTCCTCGTT
CTCAGTGTCTGTTACAGCTCGGTATTTTAGTGTCAGGACGTC
GGCTCCCAGCCCGCATCTCCGCTCAGCAATGCCATTATCTTC
TCAGCCAAGTCCTAGAAATGGGTTG**GCTTCCCATTTGCAA**
**AAACATCGCTCCATAGTCACACCAAGTCTCCTGTGGG**
**AGTCTTGCTCCTCGTTCTCAGTGTCTGTTACAGCTCG**
**GTATTTTAGTGTCAGGACGTCGGCTCCCAGCCCGCAT**
**CTCCGCTCAGCAATGCCATTATC**TTCTCAGCCAAGTCCT
AGAAATGGGTTGGCTTCCCATTTGCAAAAACATCGCTCCATA
GTCACACCAAGTCTCCTGTGGG

# Systems Biology

**DNA**
Gene 1 | Gene 2

**Transcription**

**mRNA**

CACUUCCCGAUAACGCCCCAUUGUCCUUUGGGACU
CCACGACUAUCCAGUUUCUACCUCUUCCGCAUCCCU
AUACACAACCGUCUCCUUGGCUUCUUAUGGUCCGG
UAAAGGCUUCGGUGAGUAAAAGGAGCUUAAGGAAG
AAAUA

**Translation**

**protein**

ASDVAHHPHPHAGKTALHKFTRVSLSPSNRPLLFGQT
ALHKRTFLS

XVAAAMLLRSCPVLSQGPTGLLGKVAKTYQFIFSIGRC
PILATQGPTCS

UNIVERSITEIT GENT

KU LEUVEN

# Systems Biology

# Systems Biology



Regulatory Proteins

Structural proteins

Static network is encoded in the genome

Not all proteins are made at all times (network is condition-dependent)

State of the network can be measured through gene/protein expression

# Systems Biology

computer cluster      cell culture robot      HPLC      PCR

microarray platform      MALDI-TOF mass spectrometer      DNA sequencers

# Encoding of the network

## Genome

AGCACTGTCCACTGCATGGTGAGGATGGGGGTGAGCTCCCT
TTGTGGCTAGGTGCTTAAACGTCTATCGGACGCTCA**GTGAA**
**GGGCTATTGCGGGGTAACAGGAAACCCTGAGGTGCT**
**GATAGGTCAAAGATGGAGAAGGCGTAGGGATATGTG**
**TTGGCAGAGGAACCGAAGAATACCAGGCCATTTCCG**
**AAGCCACTCATTTTCCTCGAATTCCTTCTTTAT**GCCTTC
AGTCTCTATTGACCGTAAATTTGGTTGTTGTCTCCCAGCTGT
TTATTTCTGTAACAGATCTTGGAGGCTGCGGTCTGGATCCCT
CGCCAAGAACCAGATCCAGGAGAAACGTGCTCAACGTGCA
GCTCTGCTCCTACTGATTATAGCCCCACAGATGACATCGCTC
CATAGTCACACCAAGTCTCCTGTGGGAGTCTTGCTCCTCGTT
CTCAGTGTCTGTTACAGCTCGGTATTTTAGTGTCAGGACGTC
GGCTCCCAGCCCGCATCTCCGCTCAGCAATGCCATTATCTTC
TCAGCCAAGTCCTAGAAATGGGTTG**GCTTCCCATTTGCAA**
**AAACATCGCTCCATAGTCACACCAAGTCTCCTGTGGG**
**AGTCTTGCTCCTCGTTCTCAGTGTCTGTTACAGCTCG**
**GTATTTTAGTGTCAGGACGTCGGCTCCCAGCCCGCAT**
**CTCCGCTCAGCAATGCCATTATC**TTCTCAGCCAAGTCCT
AGAAATGGGTTGGCTTCCCATTTGCAAAAACATCGCTCCATA
GTCACACCAAGTCTCCTGTGGG

# Encoding of the network

## Human Genome Project 16/02/2001



**Sequencing human genome: 13 years/ 3 miljard dollar**

# Encoding of the network



Next generation sequencing follows
Moore Law

454
Solid: 50 Gb/run;
Helicos
Illumina: 25 Gb/run; 75 bp reads

# Encoding of the network

**Next generation sequencing technology**

- **Sequencing human genome: 13 years/ 3 miljard dollar**

- **Genome Watson (454 techn): 20 persons/2 months. Totaal 1.000.000 dollar**

- **Now: 1000 dollar human genome**

# State of the network

**Functional data: transcriptome, proteome, metabolome**

**Physical data:**



Signaling network

Protein interaction network

(Post)Transcriptional network

Metabolic network

PHYSICAL INTERACTION NETWORKS

# Inferring the network



**High throughput data**



**Mechanistic insight in the biological system at molecular biological level (holistic insight)**

# Inferring the network

- **Omics data are noisy**
- **Omics data are incomplete**

- **Integrate different data to obtain higher precision and coverage**

- **Reconstruct network**
  - **Different datasources**
  - **Different Molecular layers**

# Fundamental knowledge

**Evolution: comparing network between species or over time**



Time

# Using the network to interpret data

**In house data**

**Physical interaction network**

Expression
profiling

Genes affected by
the perturbation

Infer the **hidden paths** between
genomic variations and expression
alteration

**Network-based analysis of unstructured gene lists**

UNIVERSITEIT
GENT

KU LEUVEN

# Bioinformatics and datamining

# Bioinformatics and datamining

**What does it require applying a computer science framework to bioinformatics?**

# Fast solution

- **The wet lab scientist rules**
- **Competition is fierce**
  - **Often the high impact papers are the conceptual ones**
  - **You tackle a research problem for the FIRST time**
- **Biological message is more imprtant than the method used to analyse the data**
  - **Code is sloppy , undocumented**
  - **Unsustainable code /tool development**

# Problems are complex

**Cancer subtyping & biomarker identification**



**Classification problem**

# Problems are complex



**DNA** — Gene 1 — Gene 2

Transcription

**mRNA**

**(bi)clustering problem**

# Problems are complex



genomic loci

genotypes (link data matrices)

EXPRESSION

gene

**Dataintegration problem**

# Problems are complex

## Cancer subtyping

**Preprocess the data**

genomic loci

**Associate genotypes to phenotype**

EXPRESSION

gene

**Dataintegration problem**

**Infer the interaction network from the omics data**

**Integrate all data with the interaction network**

# Bioinformatics and datamining

- **Problems need a fast solution**

- **Problems are increasingly complex and can not be solved by one particular datamining tool (generic knowledge needed)**

- **Datamining in bioinformatics requires a quite thorough understanding of biology**

- **Problems are underdetermined and ill defined**

# Network inference

**Transcriptional network inference**



**Transcription**

**Transcriptional network**

# Network inference

# Network inference

# Network inference



**Guilt by association**

**Coexpressed target genes are coregulated**

# Network inference

# Network inference

**DISTILLER**    Data Integration System    *Lemmens et al. Genome Biol. 2009*
To Identify Links in Expression in Expression regulation

## Search co expression modules that meet minimal requirements:

- All genes are significantly co-expressed in a sufficiently large, *a priori* unspecified set of experimental conditions CC

- All genes contain motif instances for a sufficient number of common, a priori unspecified regulators CR

| | A1 | A2 | A3 | … | A870 | | M1 | M2 | M3 | … | M67 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| gene 1 | 4.92 | 1.09 | -3.34 | … | 2.06 | | 0.92 | 0.86 | 0.99 | … | 0.86 |
| gene 2 | -2.21 | 0.35 | -4.98 | | 0.12 | | 0.91 | 0.95 | 0.38 | … | 0.72 |
| gene 3 | 4.45 | -2.51 | -3.98 | … | -3.41 | | 0.99 | 0.76 | 0.99 | … | 0.59 |
| … | | | | | | | | | | | |
| gene k | -4.56 | -0.13 | 4.29 | … | 2.05 | | 0.33 | 0.98 | 0.99 | … | 0.65 |

# Network inference

- Items = genes

- Transactions = conditions, motifs

| Transactions | Items | | |
|---|---|---|---|
| M1 | G1 | **G3** | **G5** |
| M2 | **G3** | **G5** | G9 |
| M3 | **G3** | **G5** | G11 |
| M4 | G1 | **G3** | **G5** |

Itemset 1: G3, G5 supported by 4 motifs
Itemset 2: G1,G3,G5 supported by 2 motifs
Tidset t(G3,G5) = M1, M2, M3, M4
Tidset t(G1, G3,G5) = M1, M4
Minimal support = 3
Itemset G3, G5 is frequent

- Supports:
  - All genes in the module (itemsets) are significantly co-expressed in a sufficiently large **(minimum support)**, a priori unspecified set of experimental conditions
  - All genes in the module (itemsets) contain motif instances for a sufficient number $R$ **(minimum support)** of common, a priori unspecified regulators

# Network inference

Expression support

co-expression in a significant number of experimental conditions

- – BW for each condition
- – Rank BW in increasing order
- – Check if BW sequence is within threshold BW sequence
- – BW threshold sequence is determined by randomization

# Network inference

Rank modules by assigning interest score

- p(module motif content):

  The chance that a module with the same number of genes and the same number of motifs is found at random

- p(module expression pattern):

  The chance that a module with at least the same number of genes and containing the same number of conditions is found by chance

  Interest score:
  p(M|m,g)Xp(M|e,g)          $\Longrightarrow$          Rank the modules

- Modules are selected iteratively such that they add as much as possible new information to the already selected modules

# CRM detection



**CRM: cis acting regulatory module**
**Combination of TF binding sites**

# CRM detection

Genes that are needed together in the cell usually are activated together = coregulation

# CRM detection

**Input:**
1) TRANSFAC PWM library
2) Set of sequences

Motif screening and filtering

(C) PWM matrix

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|------|------|------|------|------|------|
| A | 0.90 | 1.71 | -2.94 | -1.06 | -2.94 | 1.36 |
| C | -2.94 | -2.94 | -2.94 | -1.06 | 1.71 | -0.28 |
| G | 0.61 | -2.94 | -2.94 | -2.94 | -2.94 | -1.06 |
| T | -1.06 | -1.06 | 1.85 | 1.54 | -1.06 | -2.94 |

(D) sequence logo

- CRM is valid if motifs occur in each others neighbourhood
- Order needs to be conserved?

# CRM detection



DISTILLER

CPMODULE

Convert screening results in table with (start, stop) positions, for every sequence/motif pair

# CRM detection

## Constraint Programming (CP)

**Model (by user):**

Problem specification in terms of constraints

**Search (by solver):**

Propagation: in which a constraint is used to remove values from the domain of variables that would violate it

Branching: in which a variable is assigned a value from its domain D(v)



De Raedt et al. 2008 *KDD*

# CRM detection

## Frequency constraint

*Frequency constraint = 2*



Invalid CRM

Valid CRM

The motif set should occur in a sufficient number of sequences (but not all) to be considered valid

(support in itemset mining)

# CRM detection



**Proximity constraint**

50 bp

40 bp

55 bp

50 bp

40 bp

55 bp

… …

Invalid CRM

Valid CRM

Only motif instances that occur in each others proximity can contribute to a valid motif set (CRM)

# CRM detection

## Redundancy constraint

Valid CRM                    Invalid CRM



When enumerating all motif sets that meet the frequency and proximity constraint many subsets will occur in the same sequences and be composed of the same instances

Only the motif set with more motifs will be retained

Removing redundant motif sets (CRMs) drastically increases the computation time (closeness in itemset mining)

# CRM detection

Propagation: using a constraint to remove values from the domain of variables that would violate the constraint

Illustrate with the proximity constraint

propagation

Valid instances -> assign a value 1
Invalid instances-> assign a value 0

|  | m1 | m2 | m3 | ... | m66 |
|---|---|---|---|---|---|
| seq 1 | 0 | 0 | 1 | ... | 0 |
| seq 2 | 0 | 1 | 0 | ... | 0 |
| seq 3 | 1 | 0 | 1 | ... | 0 |
| ... | | | | | |
| seq k | 0 | 1 | 1 | ... | 0 |

Binary matrix dynamically updated

*Whether motif m_i is in the proximity of the motifs in motifset on sequence j?*

# CRM detection

Benchmarked on a synthetic dataset

- Xie et al. 2008 (22 sequences)
- 516 TRANSFAC PWMs
- Motifs inserted from 3 known PWMs

CPModule: performances similar to state-of-the-art algorithms on a synthetic dataset

BUT
- Able to deal with much larger sequence sets
- Enumerating all solutions allows it is able to rank the true solution amongst all solution



*Guns et al. 2010 BIBM*
*Sun et al. NAR, 2011*

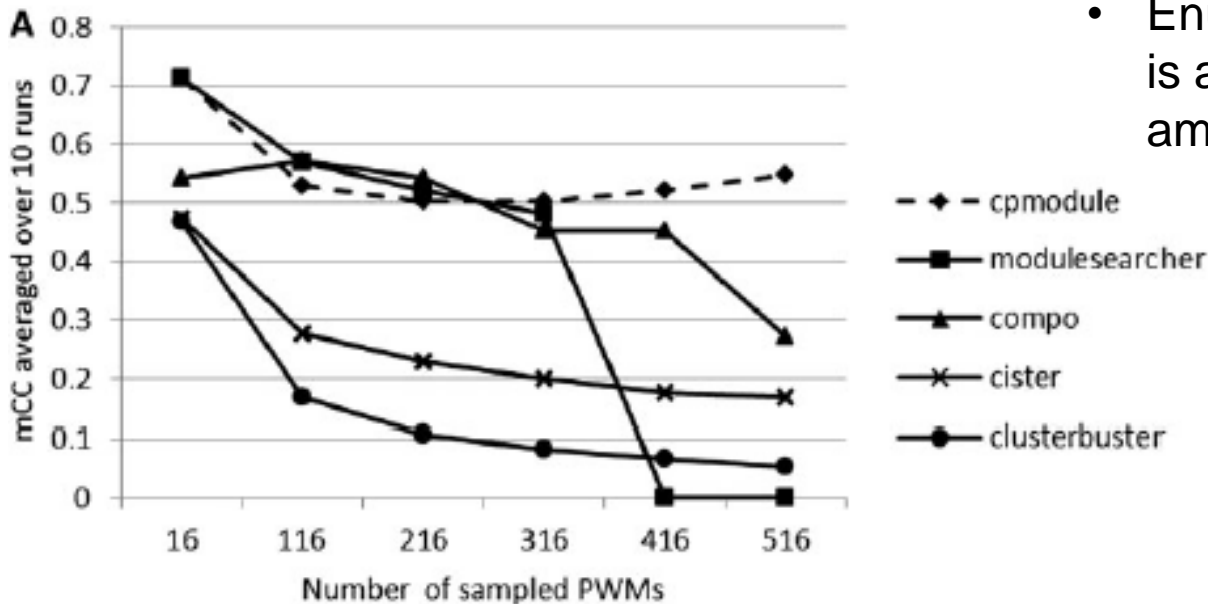# Bioinformatics and datamining

**Optimal bioinformatics tool**

- – **Right heuristics**
- – **Proper biological assumptions**
- – **Room for experimenting with different assumptions**

- – **Modular code**
- – **Sustainable code**
- – **High performance**
- – **Latest algorithmic developments**

**Usefulness of declarative framework (Problog, ASP, Constrained based programming)**

# Acknowledgements

**KUL/CMPG**

- Ivan Ischukov*
- Hong Sun*
- Valerie Storms*
- Pieter Meysman*
- Kristof Engelen*
- Lore Cloots*
- Peyman Zarrineh*
- Riet De Smet*
- Karen Lemmens*
- Abeer Fadda*

**UGENT/KUL/DBN**

- Carolina Fierro
- Yan Wu
- Aminael Sanchez
- Marleen Claeys
- Dries De Maeyer
- Sergio Pullido
- Qiang Fu

**University of Bristol**

- Tijl Debie

**KUL/Computer science**

- Luc De Raedt
- Siegfried Nijssens
- Joris Renkens
- Tias Guns
- Tan Levan

http://bioinformatics.psb.ugent.be/DBN/