

Advances in Bioinformatics

2002-2007

Jaak Vilo
vilo@ut.ee

<http://biit.cs.ut.ee>



Timeline 2002-2007

- 2002 July – move to Estonia (EGeen Inc; EBC)
- 2003 Start teaching at UT
- 2004/09 Docent (0.5) at university
- 2006/01 Senior Researcher
- 2007/12 Professor
- BIIT has now (Jan 2008)
 - 1 postdoc
 - 1 guest lecturer
 - 10 PhD students
 - 10+ students

Year	BSc	MSc
2003	2	0
2004	6	1
2005	9	1
2006	4	3
2007	4	3
Total	25	8

Research Focus



- Algorithms (Data Mining & Bioinformatics)
- Tools (web based)
- Databases & information systems

- Gene regulation & Systems Biology
- Cancer; Stem Cells;
- Microarray & other high-throughput data

Grants and Projects Completed

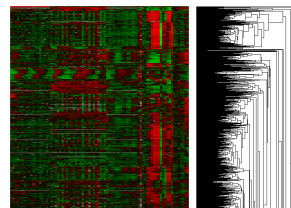
- **BiGeR** [Estonian Science Foundation](#) ETF5724 (2003-2007) Estonian Biocentre
Eesti Biokeskus, Riia 23b, Tartu 51010, Estonia
- **DMMA** [Estonian Science Foundation](#) ETF5722 (2003-2006)
- "Base funding" (start-up grant for a new research group) from University of Tartu (2005-2006)
- **ATD, Alternative Transcript Diversity**, EU FP6 STREP (2004-2007) LSHG-CT-2003-503329 Estonian Biocentre
Eesti Biokeskus, Riia 23b, Tartu 51010, Estonia
- **FunGenES, Functional Genomics of Embryonic Stem Cells**: EU FP6 Integrated Project, subcontractor (2006-2007). LSHG-CT-2003-503494 QureTEC
- [Baltic GRID](#) non-funded partner

Current grants and projects

- [COBRED](#) -- **Colon and Breast cancer Diagnostics** (2007-2010) EU FP6 STREP, LSHB-CT-2007-037730
- [ENFIN](#), **Enabling Systems Biology**. EU FP6 Network of Excellence (2005-2010) LSHG-CT-2005-518254 
- [ESNATS](#), Embryonic Stem Cell-Based Alternative Testing Strategies. Duration of Project: 5y (2008-2013). Currently under negotiation. 
- Target funding (2006-2011), The methods, environments, and applications for solving large and complex computational problems. ([SF0182712s06](#))
- ETF7437 (2008-2011) Multi-experiment gene expression data matrix analysis (MEM) ([Abstract](#))
- Estonian Language Technology Research Programme: Dictionary informatics (information retrieval) (2005-2008)

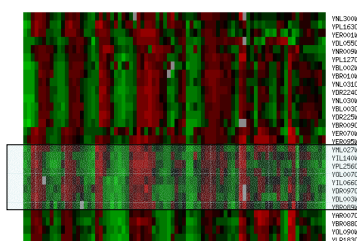
Fast Approximate Hierarchical Clustering using Similarity Heuristics

Hierarchical clustering is applied in gene expression data analysis, number of genes can be 20000+



Hierarchical clustering:

Each subtree is a cluster.

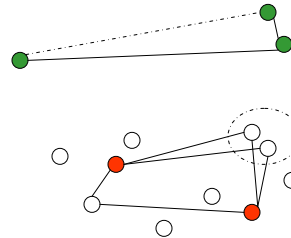


Hierarchy is built by iteratively joining two most similar clusters into a larger one.

Fast Hierarchical Clustering

Avoid calculating all $O(n^2)$ distances:

- Estimate distances
- Use pivots
- Find close objects
- Cluster with partial information



Meelis Kull

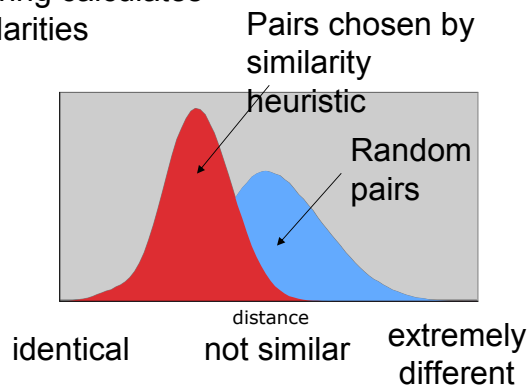


Approximate Hierarchical Clustering

- Full hierarchical clustering calculates similarities of each object with every other object (SLOW)

- Our approximate clustering calculates a small fraction of all similarities

- We have developed heuristics to find pairs of similar objects efficiently



Biological relevance

- Enrichment of (any) biological function in a cluster
- Find all enrichments in full clustering, compare quality

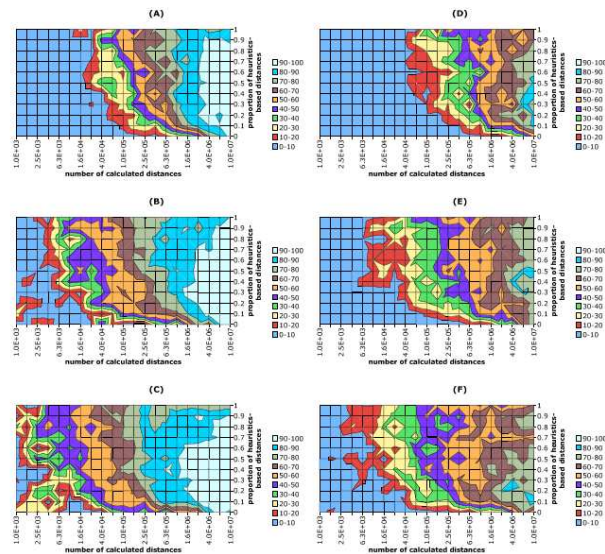
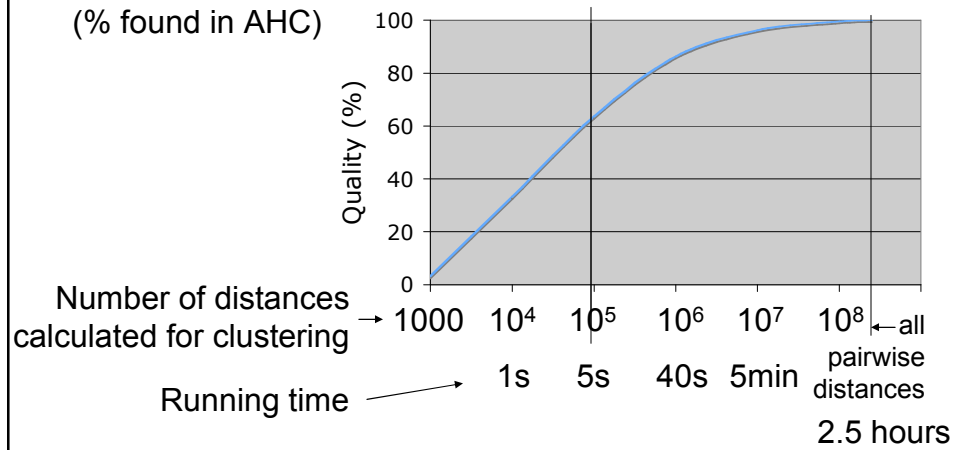
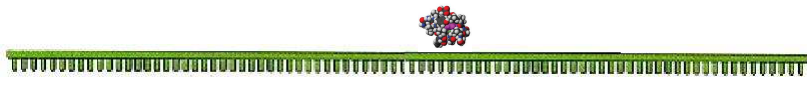
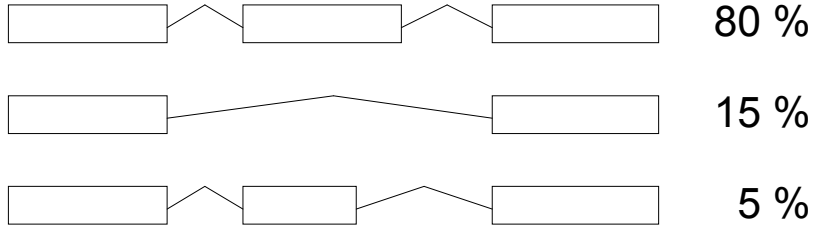
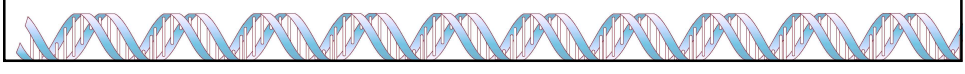


Fig. 4. GO50 and GO25 quality of HappieClust for different numbers of pivots. (A) $q = 5$, GO50; (B) $q = 10$, GO50; (C) $q = 20$, GO50; (D) $q = 5$, GO25; (E) $q = 10$, GO25; (F) $q = 20$, GO25.

Find disease-specific gene variants



Proteins can affect splicing by taking part in or interacting with the splicing complex



Results - Mozilla Firefox
 http://www.bionf.ebc.ee/~kul/atd3/pathological/

P-value: **0.00845316**
 Change: **14-fold**

Showing all results
 SHOW ONE RESULT FROM EACH GENE
 SHOW ALL RESULTS WITH GENE ENS:G00000137104

Gene name: (GeneCard)
 Ensembl ID: ENS:G00000137104 (GeneEMBL, AltSplice)
 Region: 4068..4195

Map:
 normal:
 pathological:

Zoomed map:
 normal:
 pathological:

Code	Description
Anatomical	a
Developmental	d
Pathological	p_6 Pathology---->normal

Number of ESTs	normal	pathological
4068..4195 intronic:	25 (22-26)	16 (14-17)
4068..4195 exonic:	11 (9-12)	0 (0-1)
All ESTs:	23504	32654
4068..4195 intronic per million ESTs:	106 (93-110)	48 (42-52)
4068..4195 exonic per million ESTs:	46 (38-50)	0 (0-3)

P-value: **0.0115919**
 Change: **12.4795587336096-fold**

Showing all results
 SHOW ONE RESULT FROM EACH GENE
 SHOW ALL RESULTS WITH GENE ENS:G00000137104

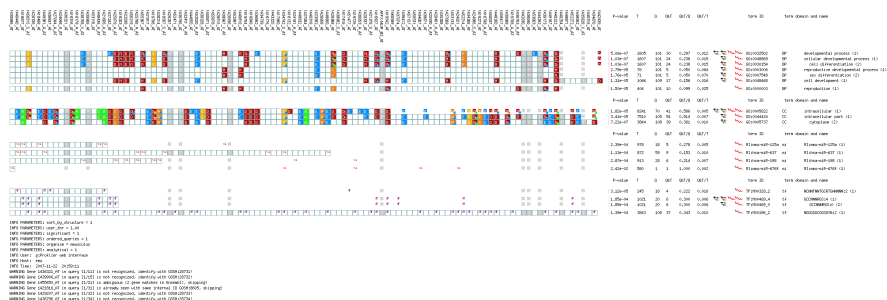
Find: splice Next Previous Highlight all Match case

http://www.ebi.ac.uk/asn-srv/atspicedb.cgi?ensembl_id=ENSG00000137104&method=ENSEMBL&specie=H&product=ALT&release=2

Meelis Kull

Tools: g:Profiler

- Characterise gene lists
- Manipulate gene ID-s
 - synonyms
 - orthologs
- Fast, rich GUI, public easy free access
- Machine-readable outputs (for integration)



Nucleic Acids Research, 2007, 1–8
doi:10.1093/nar/gkm226

g:Profiler—a web-based toolset for functional profiling of gene lists from large-scale experiments

Jüri Reimand¹, Meelis Kull^{1,2,3}, Hedi Peterson^{2,3}, Jaanus Hansen¹ and Jaak Vilo^{1,2,3,*}

¹Institute of Computer Science, University of Tartu, Liivi 2, 50409 Tartu, Estonia, ²Estonian Biocentre, Riia 23b, 51010 Tartu, Estonia and ³Egeen, Ülikooli 6a, 51003 Tartu, Estonia

Received January 31, 2007; Revised March 22, 2007; Accepted March 28, 2007

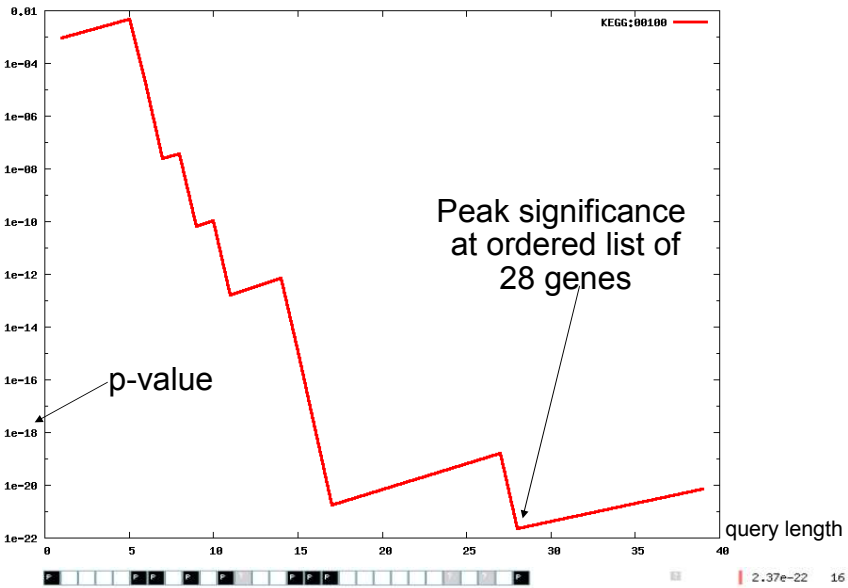
The screenshot shows the g:Profiler web application interface. At the top, there are navigation links for "g:GOST Gene Group Functional Profiling", "g:Cocoe Compact Compare of Annotations", "g:Convert Gene ID Converter", "g:Sorter Expression Similarity Search", and "g:Orth Orthology search". Below this is a "Welcome!" section with "About" and "Contact" links. The main interface is divided into several sections:

- Organism:** Homo sapiens
- Query (genes, proteins, probes):** A list of gene IDs including 210511_s_at, 212489_at, 202310_s_at, 221731_s_at, 222288_at, 208782_at, 221729_at, 61734_at, and 212344_at.
- Output options:** Includes checkboxes for "Significant only", "Hierarchical sorting", "1e-5 User p-value", and "Output type" (set to Graphical (PNG)).
- Input options:** Includes checkboxes for "Ordered query", "Ignore unknown entries", "Show advanced options", "Direct assay [IDA] / Mutant phenotype [IMP]", "Genetic interaction [IGI] / physical interaction [IP1]", "Expression pattern [IEP] / Reviewed computational analysis [RCA]", "Sequence or structural similarity [ISS] / Genomic context [IGC]", "Traceable author [TAS] / Inferred by curator [IC]", "Non-traceable author [NAS]", "Electronic annotation [IEA]", "Multiple GO evidence codes", "No data [ND] / Not annotated", "KEGG/REACTOME pathway", "TRANSFAC regulatory motifs", and "miRBase microRNAs".
- Buttons:** "g:Profile" and "Clear".
- Navigation:** "g:Convert Gene ID Converter", "g:Orth Orthology Search", "g:Sorter Expression Similarity Search", and "Static URL Come back later".
- Results Table:** A table with columns: P-value, T, O, O&T, O&T/O, O&T/T, term ID, and term domain and name. The table shows two rows of results with various GO terms like "ion transport (1)", "anion transport (2)", "inorganic anion transport (3)", "phosphate transport (4)", "anatomical structure development (1)", "system development (2)", "organ development (2)", and "skeletal development (3)".

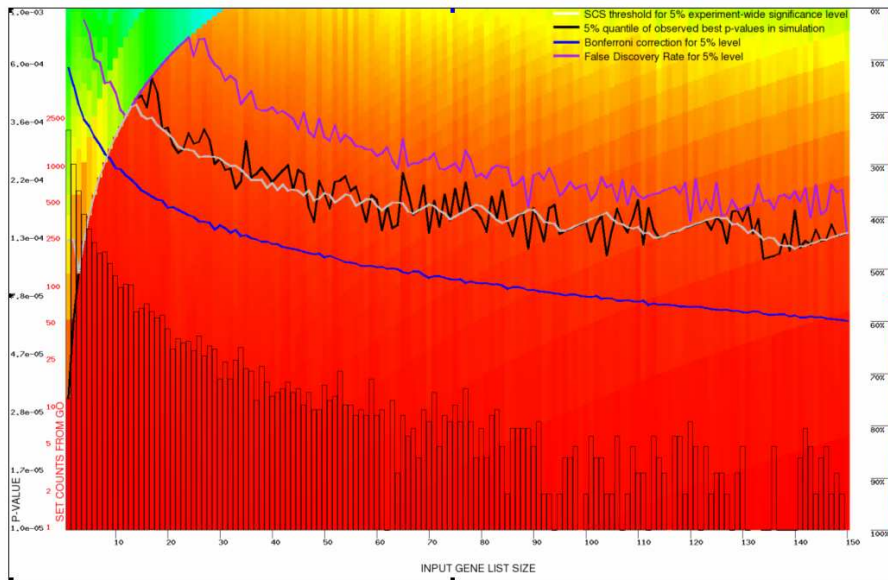
This screenshot is similar to the one above but includes several blue callout boxes highlighting specific features:

- Genes:** A callout box pointing to the "Query (genes, proteins, probes)" list.
- Evidence codes:** A callout box pointing to the "Evidence codes" section in the "Input options" area.
- GO: KEGG, Reactom, microRNA:** A callout box pointing to the "KEGG/REACTOME pathway", "TRANSFAC regulatory motifs", and "miRBase microRNAs" options.
- P-value:** A callout box pointing to the "P-value" column in the results table.
- Ordered list query:** A callout box pointing to the "Ordered query" checkbox in the "Input options" area.

KEGG: Biosynthesis of



SCS - Set Counts and Sizes



27	7	1.76e-11	581	BP	GO:0007155	BP	cell adhesion (1)
		4.46e-05	23	BP	GO:0006775	BP	Fat-soluble vitamin metabolic process (1)
		2.08e-05	18	BP	GO:0006776	BP	vitamin A metabolic process (2)
7	2	1.44e-05	49	BP	GO:0045997	BP	positive regulation of cell differentiation (1)
15		5.92e-06	173	BP	GO:0045892	BP	negative regulation of transcription, DNA-dependent (1)
7		4.21e-06	41	BP	GO:0043406	BP	positive regulation of MAPK activity (1)
42	13	4.60e-08	4146	BP	GO:0065907	BP	biological regulation (1)
42	13	3.56e-07	3834	BP	GO:0050789	BP	regulation of biological process (2)
14	24	4.87e-06	895	BP	GO:0048519	BP	negative regulation of biological process (3)
13	43	7.19e-15	2605	BP	GO:0032502	BP	developmental process (1)
12	34	2.48e-14	1733	BP	GO:0048856	BP	anatomical structure development (2)
9	22	3.59e-13	997	BP	GO:0009933	BP	anatomical structure morphogenesis (2)
2	3	3.95e-05	57	BP	GO:0001763	BP	morphogenesis of a branching structure (3)
2	3	2.44e-05	93	BP	GO:0048754	BP	branching morphogenesis of a tube (4)
20	11	3.94e-05	437	BP	GO:0032989	BP	cellular structure morphogenesis (3)
20	11	3.94e-05	437	BP	GO:0000902	BP	cell morphogenesis (4)
5	14	2.95e-07	185	BP	GO:0048646	BP	anatomical structure formation (2)
16		3.93e-08	232	BP	GO:0050793	BP	regulation of developmental process (2)
13	38	8.30e-14	1937	BP	GO:0007275	BP	multicellular organismal development (1)
5	11	4.54e-05	398	BP	GO:0009790	BP	embryonic development (2)
3	9	3.24e-06	189	BP	GO:0007389	BP	pattern specification process (2)
7		2.93e-05	139	BP	GO:0003002	BP	regionalization (3)
6		2.15e-05	89	BP	GO:0009952	BP	anterior/posterior pattern formation (4)
6		3.74e-06	144	BP	GO:0035295	BP	tube development (2)
11	27	1.95e-11	1491	BP	GO:0048731	BP	system development (2)
8	24	3.42e-10	1202	BP	GO:0048513	BP	organ development (3)
6	16	1.72e-08	207	BP	GO:0001944	BP	vasculature development (4)
6	16	1.43e-08	204	BP	GO:0001568	BP	blood vessel development (5)
7	18	8.25e-08	441	BP	GO:0009887	BP	organ morphogenesis (4)
5	13	2.14e-08	177	BP	GO:0048514	BP	blood vessel morphogenesis (5)
11	10	1.15e-06	139	BP	GO:0001525	BP	angiogenesis (6)
14	23	5.05e-08	1607	BP	GO:0048869	BP	cellular developmental process (1)
14	23	5.05e-08	1607	BP	GO:0003054	BP	cell differentiation (2)
14	23	3.12e-06	3499	BP	GO:0050794	BP	regulation of cellular process (1)
6	13	6.11e-06	836	BP	GO:0048523	BP	negative regulation of cellular process (2)
2	7	6.23e-07	78	BP	GO:0003006	BP	reproductive developmental process (1)
6		5.82e-06	71	BP	GO:0007548	BP	sex differentiation (2)
14	23	5.72e-06	1066	BP	GO:0048468	BP	cell development (1)
7	10	3.91e-08	226	BP	GO:0007167	BP	enzyme linked receptor protein signaling pathway (1)
12		7.64e-06	151	BP	GO:0007169	BP	transmembrane receptor protein tyrosine kinase signaling pathway (2)
24	6	7.26e-11	647	BP	GO:0002376	BP	immune system process (1)
13	6	4.08e-06	370	BP	GO:0006992	BP	defense response (1)
2	14	2.22e-07	336	BP	GO:0042221	BP	response to chemical stimulus (1)
12		1.29e-07	225	BP	GO:0001775	BP	cell activation (1)
10		3.99e-06	210	BP	GO:0045321	BP	leukocyte activation (2)
9	20	5.24e-11	425	BP	GO:0009605	BP	response to external stimulus (1)
12		1.27e-06	280	BP	GO:0007610	BP	behavior (1)
12		8.87e-09	177	BP	GO:0007626	BP	locomotory behavior (2)
4	11	2.37e-10	101	BP	GO:0042330	BP	taxis (3)
4	11	2.37e-10	101	BP	GO:0006935	BP	chemotaxis (4)
13	6	2.34e-10	412	BP	GO:0006955	BP	immune response (1)
13	6	1.98e-07	277	BP	GO:0009611	BP	response to wounding (1)

Data Analysis Environments

- FunGenES consortium – data analysis environment
- ~10 labs produced 150+ conditions (3-5x)
 - Affymetrix Gene Chips (22,000 probesets; ~30MB file each)
 - Normalise, extract data for analysis
 - study, analyse, visualise, ...
 - compare
- Clustering, visualisation, searching, etc.

FunGenES project @ Tartu - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://bit.cs.ut.ee/fungenes/

FunGenES project @ Tartu FunGenES article supplement...

FunGenES
Functional Genomics
in Embryonic
Stem Cells

g:Profiler
g:Profiler
g:Orth
g:Convert

Visualize the gene list
Heatmapper
URLMAP

KEGG pathway animations
KEGG animations

Supplementary for common publication
Database

Clusterings
Individual datasets
Global datasets

Expression waves
Similarity cutoff 0.8
Similarity cutoff 0.85

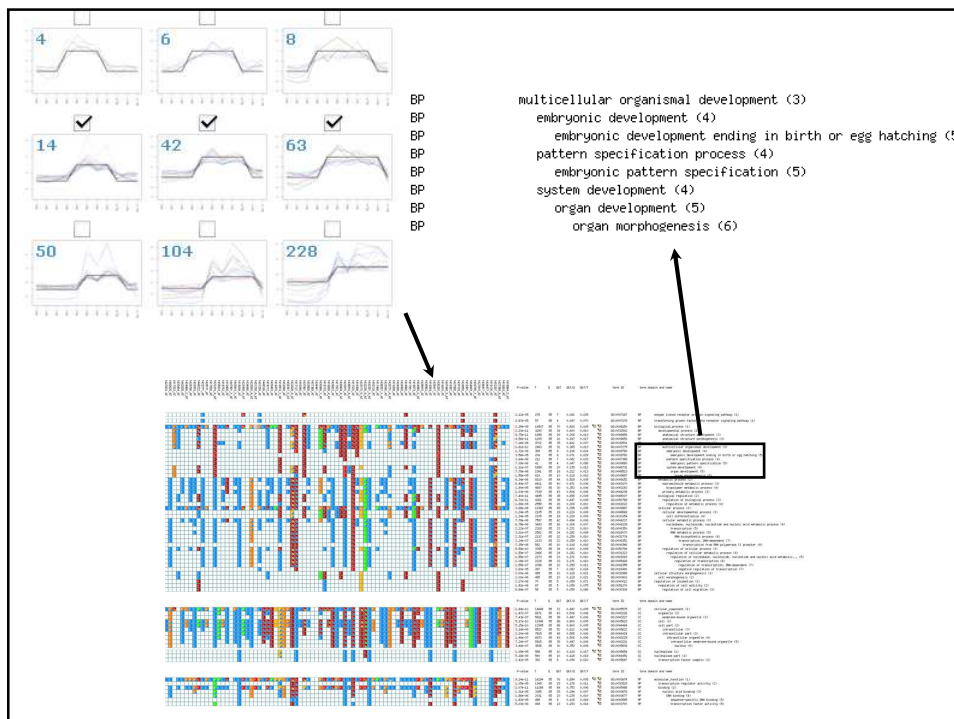
KEGG pathway animations
KEGG animations

Supplementary for the article
Database

Supporting website for common publication.

Clusterings
Clusterings of individual datasets
Clusterings of global datasets
Clusterings of the FunGenES data. Individual datasets and global combinations of those.

Find: Next Previous Highlight all



Regulators for ES gene clusters

DELETED

FunGenES

- Analysed data for 4-5 groups; others used web
- 4 people visited Estonia (for a week) to analyse their data
- Consortium data – public analysis environment
- Couple of leads in functional validation

Tools: KEGGanim

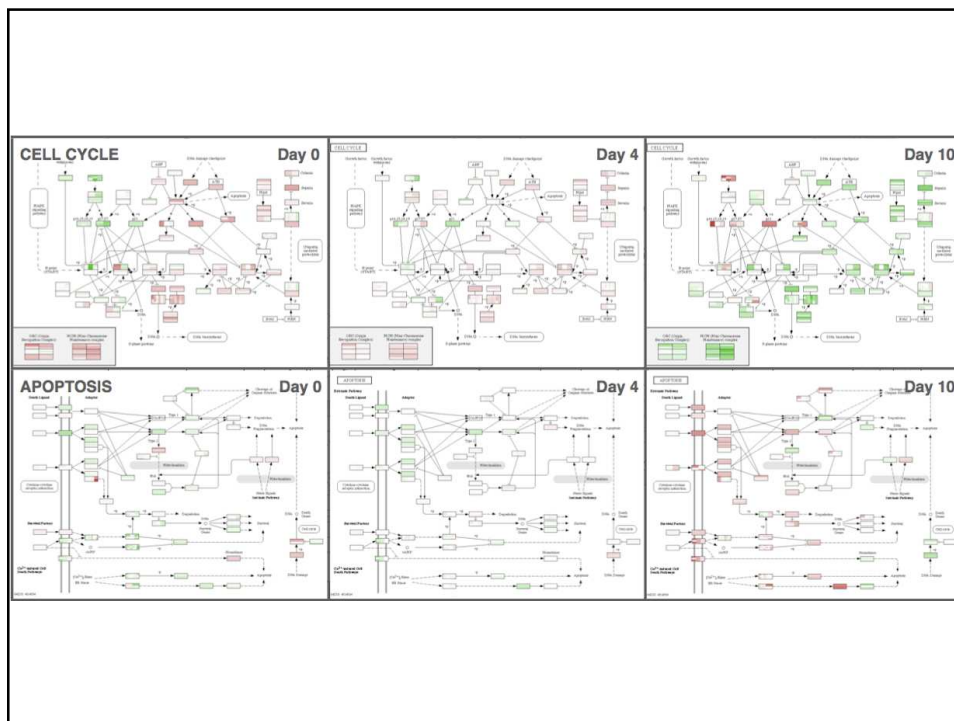
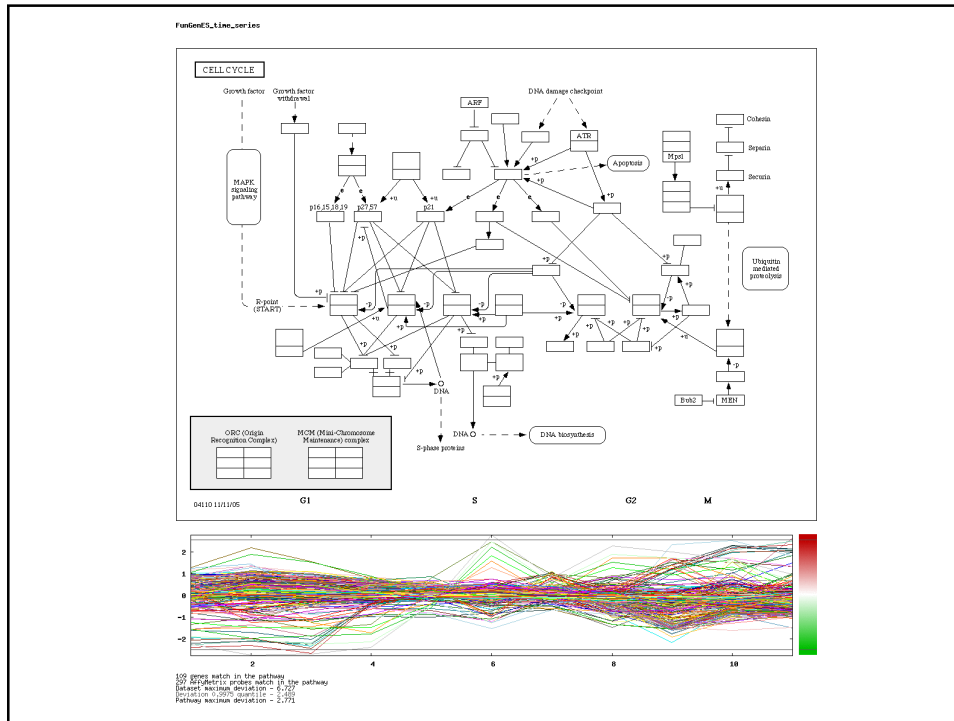
- Animate large-scale mRNA, proteomics, metabolomics data on KEGG pathways
- Create visualisations
 - animations
 - freeze panels
- Manipulate private data

biit.cs.ut.ee/KEGGanim
Adler et al
Bioinformatics, 2007

Overlay HT data over KEGG pathways

- Animate
- “Freeze panes” (cinofilm)

KEGGanim: pathway animations for high-throughput data
Prit Adler^{a*}, Jüri Reimand^{b*}, Jürgen Jänes^b, Raivo Kolde^c, Hedi Peterson^{ac}, Jaak Vilo^{abct}
^aEstonian Biocentre, Riia 23b, Tartu, Estonia ^bUniversity of Tartu, Institute of Computer Science, Livi 2, Tartu, Estonia ^cQureTec Inc, Üikooli 6a, Tartu, Estonia

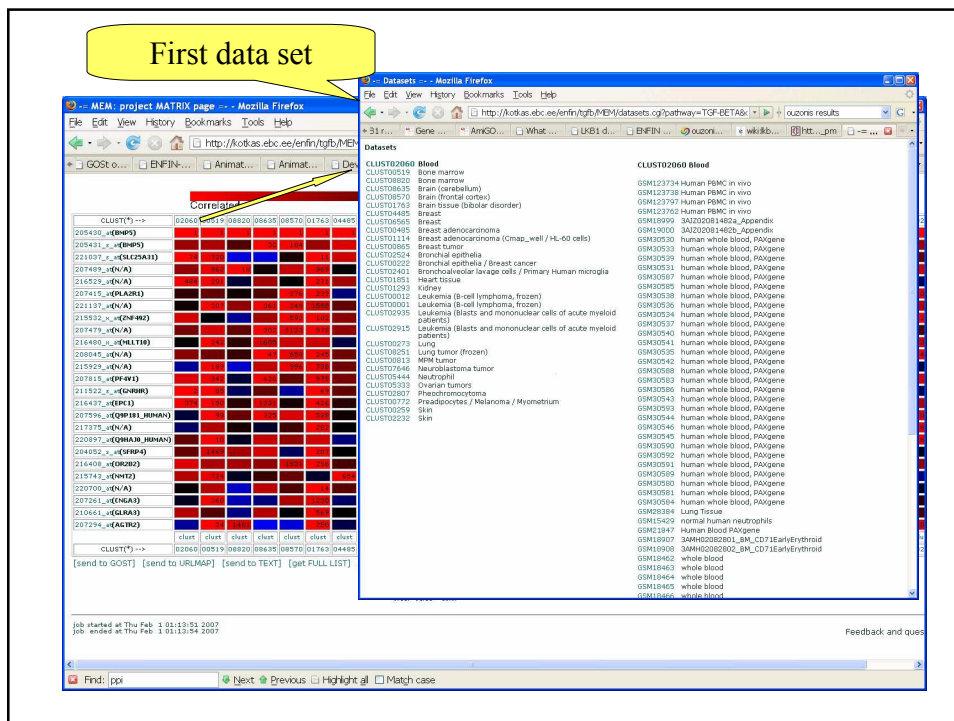
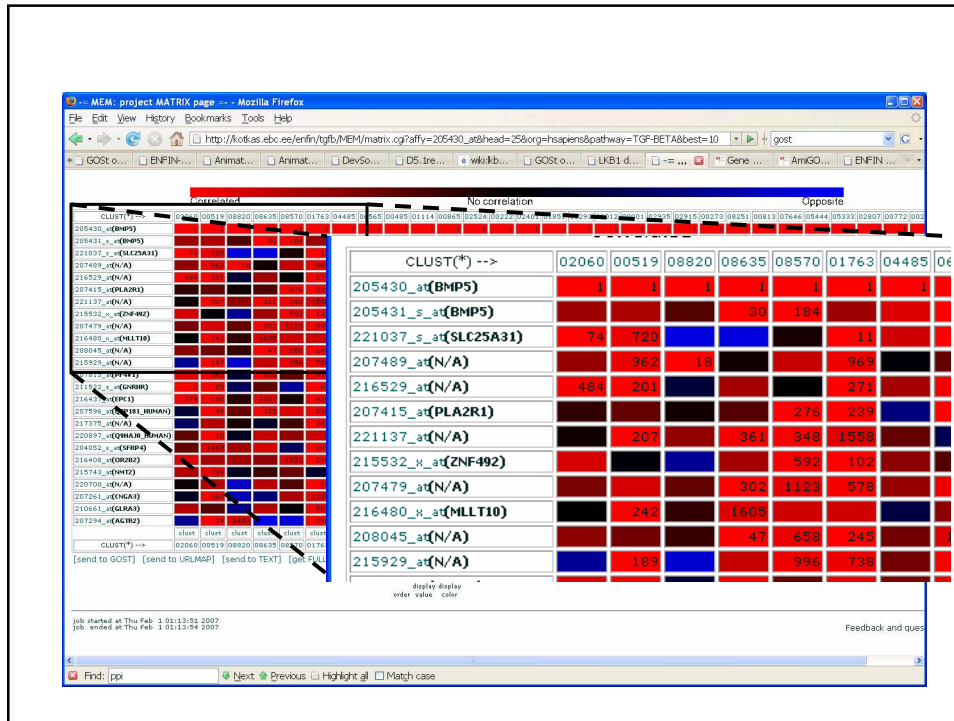


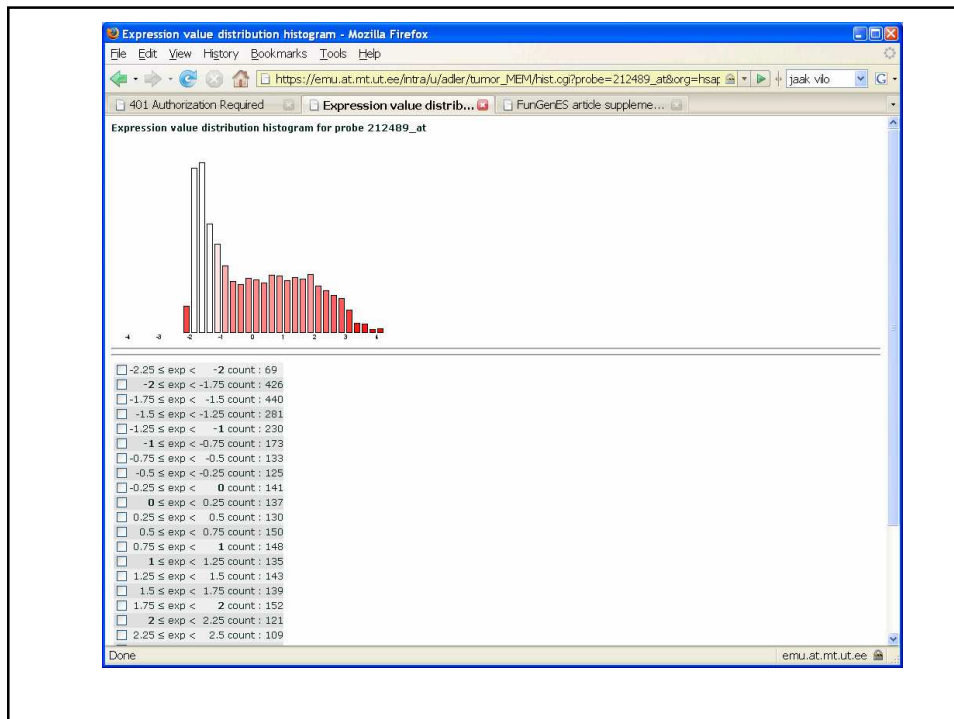
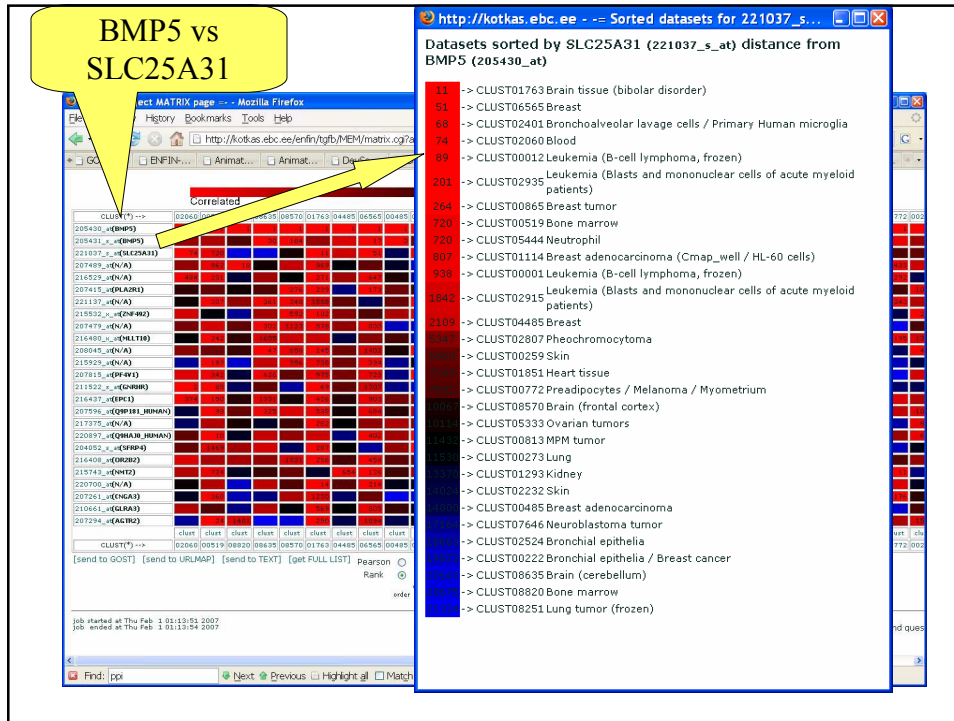
More large-scale experimental data

- Gene expression – thousands of experiments characterising 20K+ genes
- Protein-Protein interactions (complexes, signal transduction, ...)
- Protein-DNA interactions
- proteomics, metabolomics, genetics, ...

MEM – Multiple Experiment Matrix Co-expression of TGF β gene(s) across multiple datasets



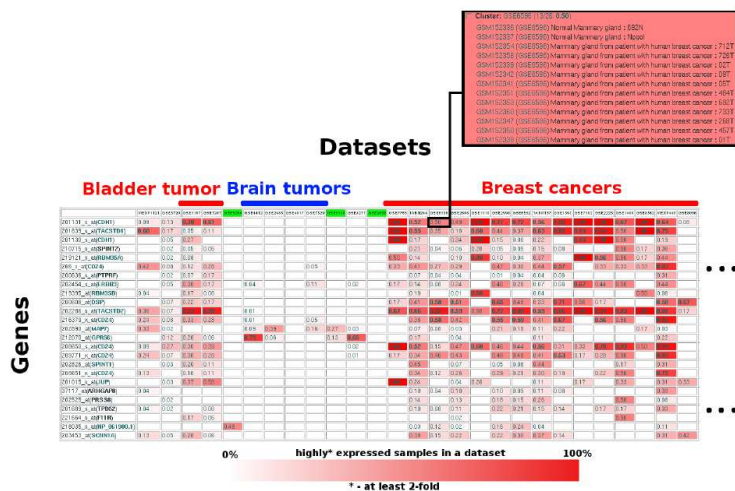




MEM goals

- Query across hundreds of datasets simultaneously
- Characterise gene expression across all datasets
- Study in detail all genes within one pathway
- Find new genes possibly linked to the given pathway
- Identify relevant “conditions” for gene activity
- Priit Adler; Raivo Kolde; Hedi Peterson

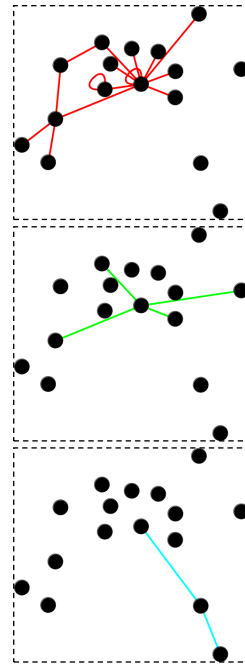
MEM - applications - global expression distribution Oncogene characterisation



Hedi Peterson MEM - global expression analysis tool

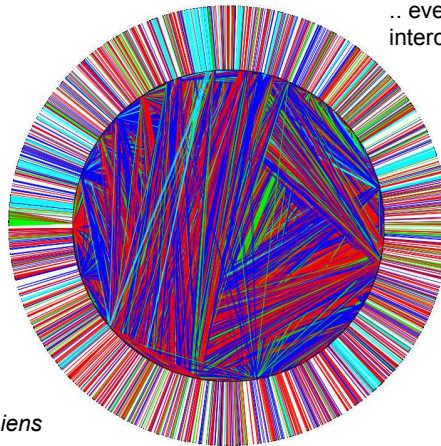
GraphWeb: mining biological networks for submodules with functional significance

- Genes as nodes
 - -omics define edges
- | | | |
|---|---|------------------------------|
| ◇ | — | expression correlation |
| ◇ | — | protein-protein interactions |
| ◇ | — | literature co-occurrence |
| ◇ | — | regulation |
| ◇ | — | binding site discovery |



Data as graphs

.. everything is interconnected

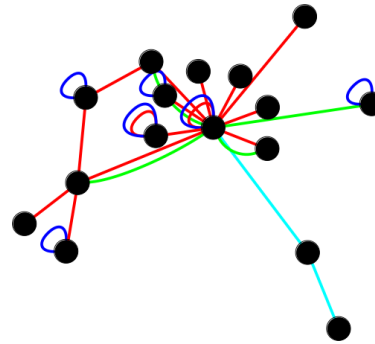


Public datasets for *H.sapiens*

- ◇ IntAct: Protein interactions (PPI), 18773 interactions
- ◇ IntAct: PPI via orthologs from IntAct, 6705 interactions
- ◇ MEM: gene expression similarity over 89 tumor datasets, 46286 interactions
- ◇ Transfac: gene regulation data, 5183 interactions

Gene modules

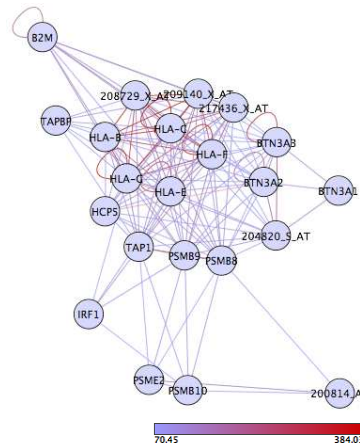
- Integrate data sources as graph layers
- Find well-connected subgraphs
- Combine evidence to infer knowledge about regulation and function



GO: cell cycle, regulation, growth.
KEGG: Alzheimer's disease

Weighing the evidence I

- Edges are not born equal
 - e.g. stronger vs weaker correlation
- Assign *local* weights to rank edges within a layer
- Look for *heavy* subgraphs



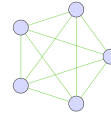
Expression similarity search across 89 human tumor-related datasets (MEM, Adler et al, in prep.)

GO: immune system, proteasome. Reactome: cell cycle, DNA replication, HIV infection

Finding the modules

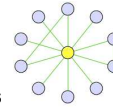
- **Cliques**

- Fully connected graphs ~ protein complexes



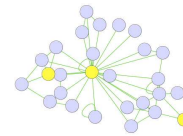
- **Hubs**

- Highly connected nodes ~ transcriptional regulators



- **Sets of neighbors**

- Specific genes of interest + near neighbors



- **Graph clustering**

- MCL: Markov clustering (*van Dongen, 2000*), betweenness centrality clustering

Module evaluation

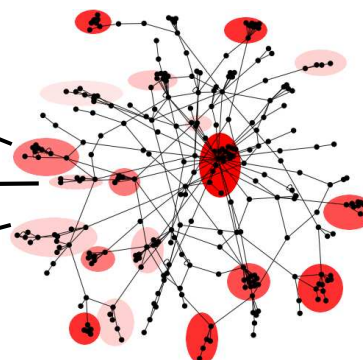
GO: Transforming growth factor beta signaling pw.
embryonic development, gastrulation
KEGG: Cell cycle, cancers, WNT pw.



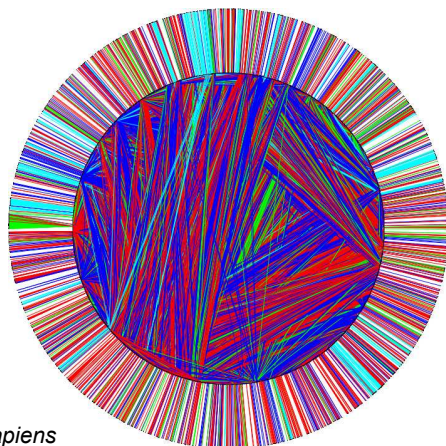
GO: JAK-STAT cascade, Kinase inhibitor activity
Insulin receptor signaling pw.
KEGG: Type II diabetes mellitus



GO: Brain development
Pigment granule
Melanine metabolic process



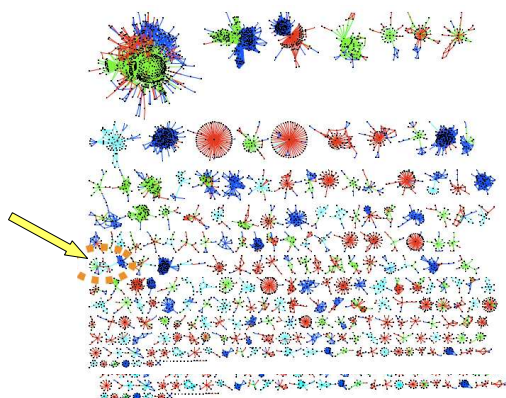
Finding the modules



Public datasets for *H.sapiens*

- ◆ IntAct: Protein interactions (PPI), 18773 interactions
- ◆ IntAct: PPI via orthologs from IntAct, 6705 interactions
- ◆ MEM: gene expression similarity over 89 tumor datasets, 46286 interactions
- ◆ Transfac: gene regulation data, 5183 interactions

Finding the modules



Public datasets for *H.sapiens*

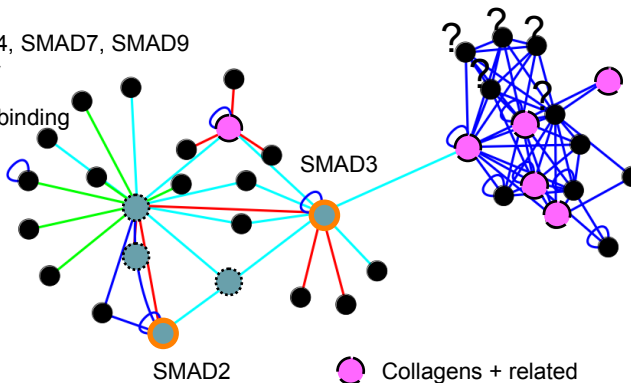
- ◆ IntAct: Protein interactions (PPI), 18773 interactions
- ◆ IntAct: PPI via orthologs from IntAct, 6705 interactions
- ◆ MEM: gene expression similarity over 89 tumor datasets, 46286 interactions
- ◆ Transfac: gene regulation data, 5183 interactions

Example

● SMAD2, SMAD3, SMAD4, SMAD7, SMAD9
transcription factor family

● SMAD2, SMAD3?: DNA binding

Transforming Growth
Factor Beta signalling,
Development, Cancers



Public datasets for *H.sapiens*

- ◆ IntAct: Protein interactions (PPI), 18773 interactions
- ◆ IntAct: PPI via orthologs from IntAct, 6705 interactions
- ◆ MEM: gene expression similarity over 89 tumor datasets, 46286 interactions
- ◆ Transfac: gene regulation data, 5183 interactions

GraphWeb

<http://biit.cs.ut.ee/graphweb>

Input: a large graph with experimental data

GraphWeb

<http://biit.cs.ut.ee/graphweb>

Module ID	Nodes	Edges	Zoom	Bar Chart	Score	Support	Annotations	Visuals	
# 18	43	List nodes List edges	Zoom in	[Bar Chart]	25	282	7	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 19	42	List nodes List edges	Zoom in	[Bar Chart]	11	145	3	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 20	42	List nodes List edges	Zoom in	[Bar Chart]	64	710	17	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 21	41	List nodes List edges	Zoom in	[Bar Chart]	1	9	0	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 22	40	List nodes List edges	Zoom in	[Bar Chart]	49	969	25	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 23	40	List nodes List edges	Zoom in	[Bar Chart]	30	579	14	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 24	38	List nodes List edges	Zoom in	[Bar Chart]	66	1776	47	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 25	38	List nodes List edges	Zoom in	[Bar Chart]	8	92	2	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 26	38	List nodes List edges	Zoom in	[Bar Chart]	0	0	0	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 27	35	List nodes List edges	Zoom in	[Bar Chart]	1	8	0	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 28	34	List nodes List edges	Zoom in	[Bar Chart]	39	532	16	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 29	34	List nodes List edges	Zoom in	[Bar Chart]	27	711	21	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 30	34	List nodes List edges	Zoom in	[Bar Chart]	1	14	0	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 31	33	List nodes List edges	Zoom in	[Bar Chart]	5	52	2	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 32	31	List nodes List edges	Zoom in	[Bar Chart]	39	730	24	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 33	31	List nodes List edges	Zoom in	[Bar Chart]	7	66	2	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 34	29	List nodes List edges	Zoom in	[Bar Chart]	1	9	0	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 35	29	List nodes List edges	Zoom in	[Bar Chart]	0	0	0	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi
# 36	27	List nodes List edges	Zoom in	[Bar Chart]	7	72	3	g:P structured flat	labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi

Output: a list of tightly connected gene modules

GraphWeb

<http://biit.cs.ut.ee/graphweb>

36 List nodes List edges Zoom in [Bar Chart] 95 4272 119 g:P structured flat labeled: dot neato fdp circo twopi compact: dot neato fdp circo twopi

Export graph nodes/edges or zoom in GraphWeb

Evaluate edge content and support from datasets

Graph scores from functional annotations

g:Profiler

Visuals

Protein interactions (PI) 8%

8% 5% 8% 2%

FPI via orthologs

MEM expression 3%

GO: catabolic process, proteasome
KEGG: DNA replication, HIV infection, Cell cycle checkpoints

microRNA - discovery

- Alexander Stark, Pouya Kheradpour, **Leopold Parts**, Julius Brennecke, Emily Hodges, Gregory J. Hannon, and Manolis Kellis. Systematic discovery and characterization of fly microRNAs using 12 *Drosophila* genomes. *Genome Research*, 2007
- Alexander Stark, Michael F. Lin, Pouya Kheradpour, Jakob S. Pedersen, **Leopold Parts**, et. al. Discovery of functional elements in 12 fly genomes using evolutionary signatures. *Nature* **450**, 219-232 (8 November 2007)
- *Drosophila* 12 Genomes Consortium. Evolution of genes and genomes on the *Drosophila* phylogeny. *Nature*, 2007

Next 5 years

- Databases and analysis of human health
 - Colon & Breast cancer; disease relapse
 - mRNA, proteomics, metabolomics data
 - Estonian Genome Project
 - DB of 10's of thousands of patients with ~1000 attributes worth of health and lifestyle profiles
 - ES cell based toxicology profiling for new drug candidates
 - Disease/gene associations

Next 5 years

- Data query and visualisation methods
 - MEM, PATMATCH, ...
- Data mining
 - Motif discovery (SPEXS, Trie*Tools, ...)
 - Gene regulation modeling
 - Phenotype, Genotype, Expression, PPI, CHIP-chip, ... data joint analysis
- Data analysis environments
- Analysis of data

Anno 2007 (BIIT and Quretec)



Institute of Computer Science

- ★ Software Engineering (Prof. Marlon Dumas)
 - High-Performance (distributed) computing (Prof. Eero Vainikko)
 - Language Technology (Prof. Mare Koit)
- ★ Programming Languages (Prof. Varmo Vene)
 - Cryptography (Prof. Ahto Buldas)
- ★ Bioinformatics and DM (Prof. Jaak Vilo)
 - 2 Target funding projects (HPC and LT)
- ★ new appointments in Nov & Dec 2007