

Gene regulation, protein networks and disease - a computational perspective

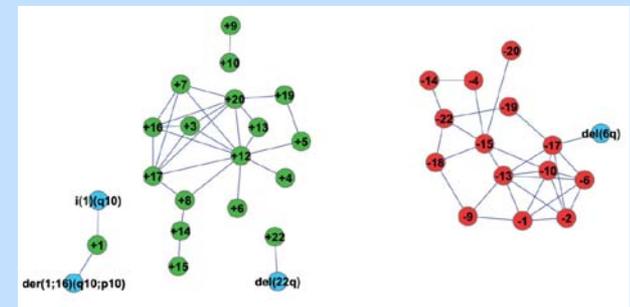
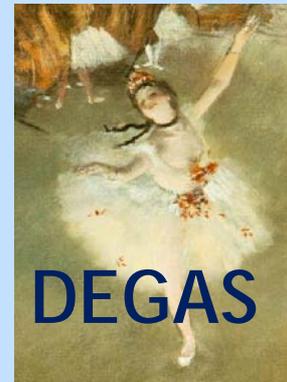
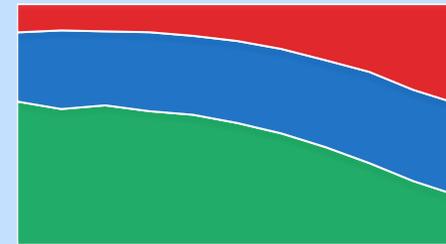
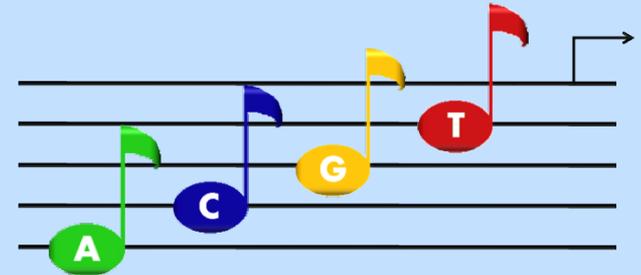
Ron Shamir
School of Computer Science
Tel Aviv University

CPM Helsinki July 3 2012



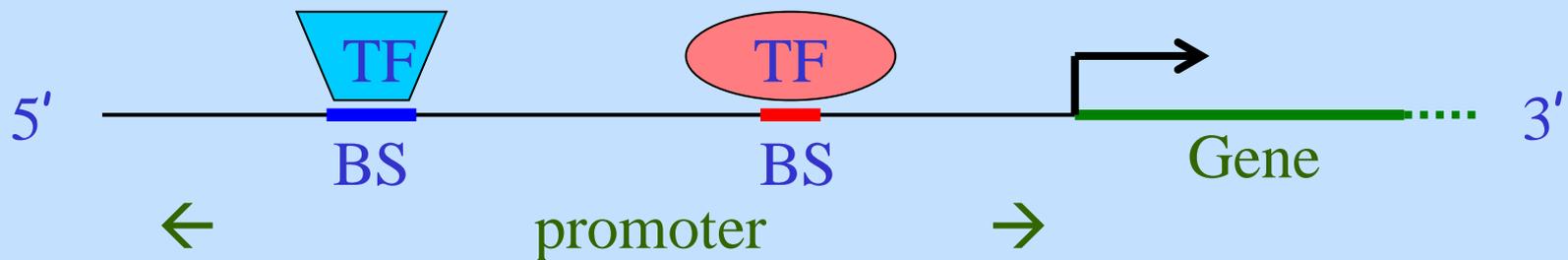
Outline

- Finding regulatory motifs I, II, III
- Utilizing case-control expression profiles and networks I, II
- Chromosomal aberrations in cancer



Regulation of Transcription

- A gene's transcription regulation is mainly encoded in the DNA in a region called the **promoter**
- Each promoter contains several short DNA subsequences, called **binding sites (BSs)** that are bound by specific proteins called **transcription factors (TFs)**



Position Weight Matrix (PWM)

Score: product of base probabilities.

Need score threshold for hits.

A	0.1	0.8	0	0.7	0.2	0
C	0	0.1	0.5	0.1	0.4	0.6
G	0	0	0.5	0.1	0.4	0.1
T	0.9	0.1	0	0.1	0	0.3

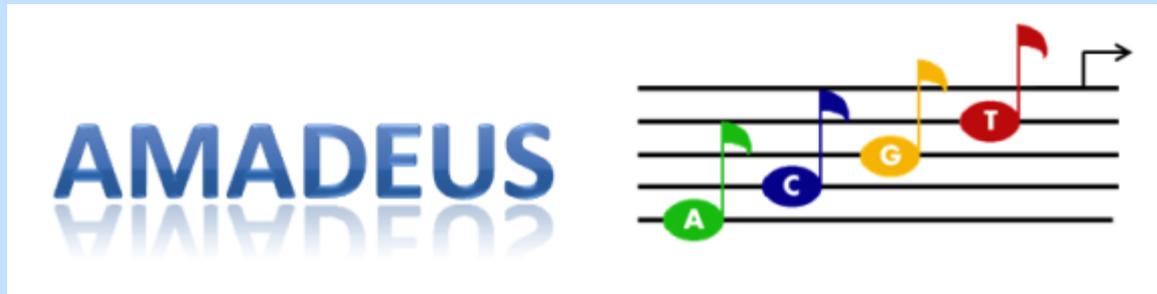
ATGCAGGAT**TACACCG**ATCGGTA 0.0605

GGAG**TAGAGCA**AGTCCCGTGA 0.0605

AAGACTC**TACAAT**TATGGCGT 0.0151



I. Finding Regulatory Motifs

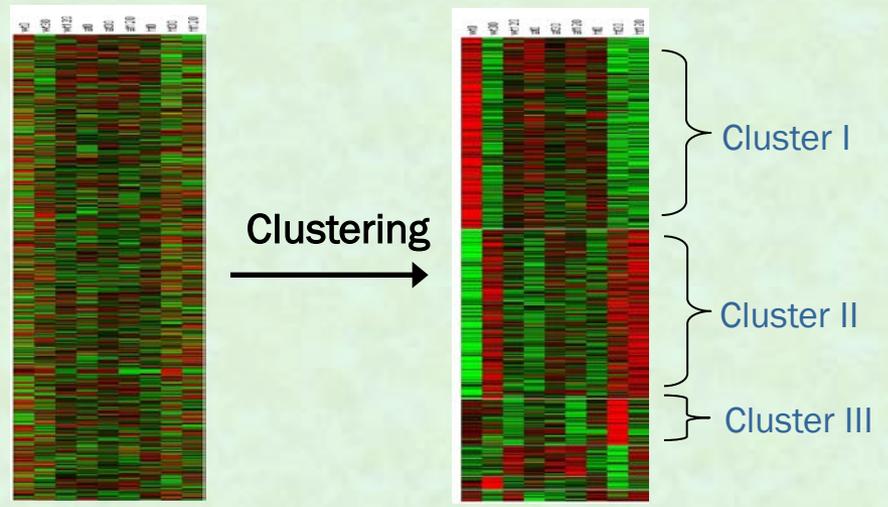


Motif discovery:

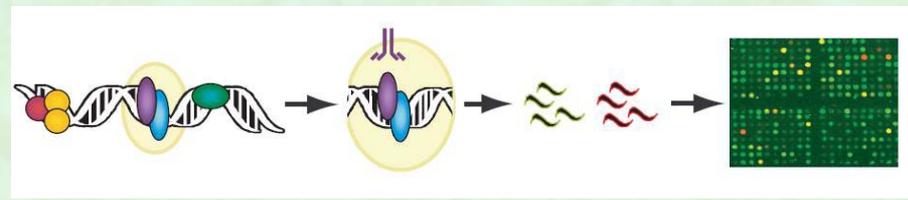
The two-step strategy

Co-regulated gene set

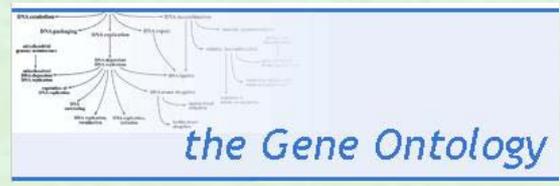
Gene expression
microarrays



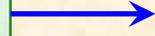
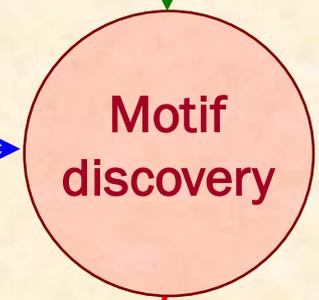
Location analysis
(ChIP-chip, ...)



Functional group
(e.g., GO term)



Promoter
sequences





Amadeus

A Motif Algorithm for Detecting Enrichment in Multiple Species

- Supports diverse motif discovery tasks:
 1. Find **over-represented** motifs in given **sets of genes**.
 2. Identify motifs with **global spatial features** given **only** the genomic **sequences**.
- **How?**
 - A general **pipeline architecture** for enumerating motifs.
 - Different statistical **scoring schemes** of motifs for different motif discovery tasks.



Motif search algorithm

- Pipeline of refinement phases of increased complexity

➤ Phases:



CGGACTTTCC
GACTTAGACT
CGGATCGATT
GATAGTACCG
CCATATCCGA
AATTTGAATC
AACCCGTGCA
TGCGATTACG

GATAGTACCG
↓
G_eTAGTAcCG

GAAATCC
+
G₁AA₂T₁TC₁
↓
G₁AA₂T₁TC₁

EM
GG₁TTCC
Cutoff = 0.005

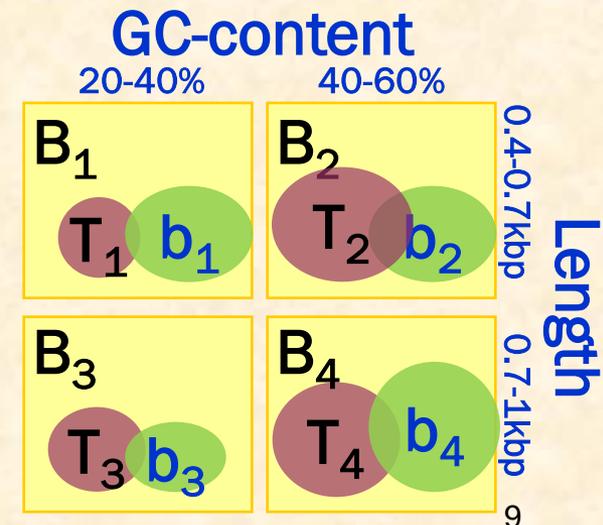
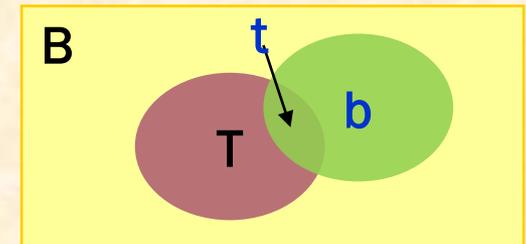


Motif Model:



Scoring over-represented motifs

- Input: Target set (size T) = co-regulated genes
 Background (BG) set (size B) = entire genome
- Motif enrichment scoring:
- Hyper-geometric
- Binned enrichment score
- Binomial



Metazoan motif discovery benchmark:

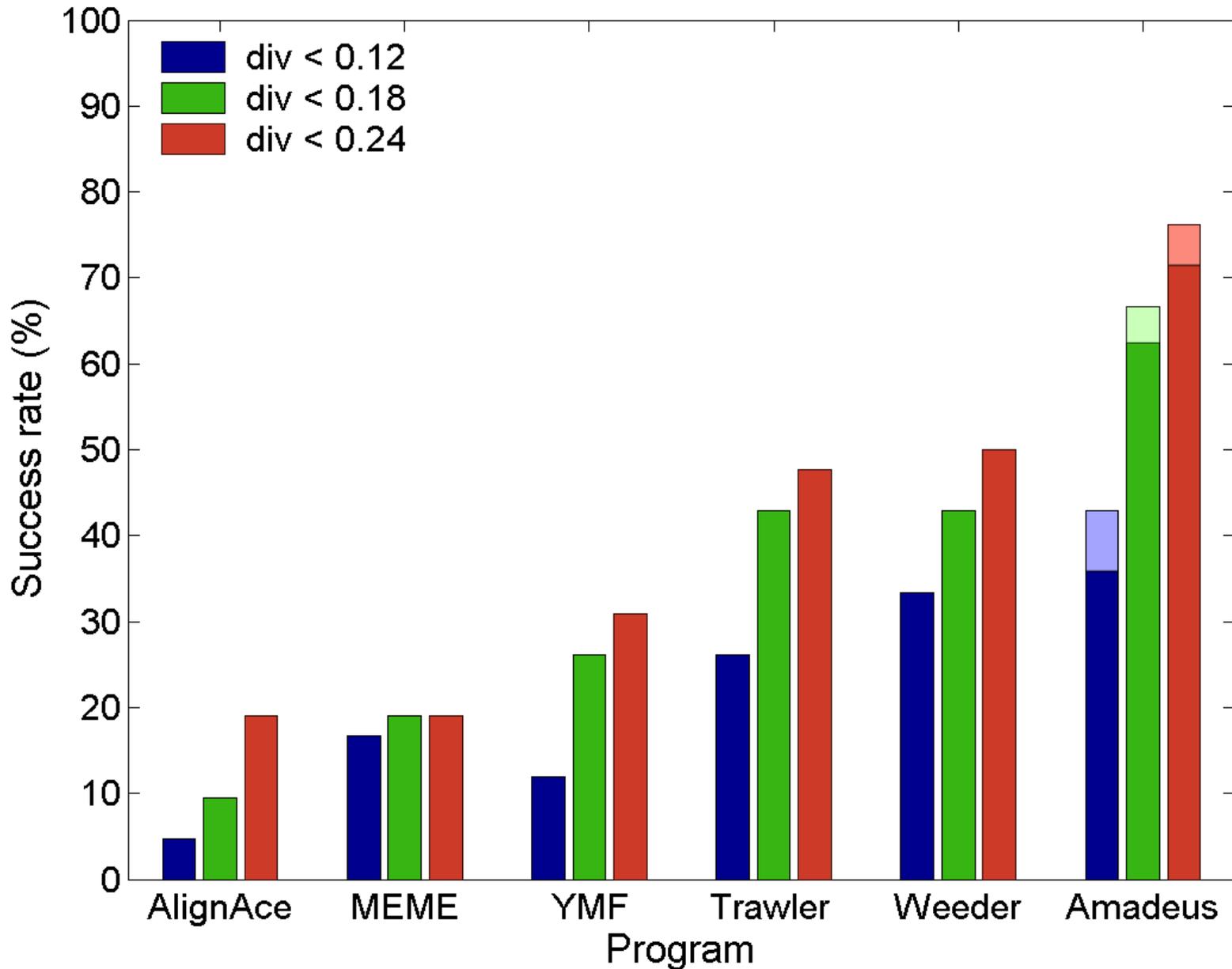
42 target sets of 26 TFs, 8 miRNAs from 29 studies
(expression, Chip-ChIP,..) in human, mouse, fly, worm.

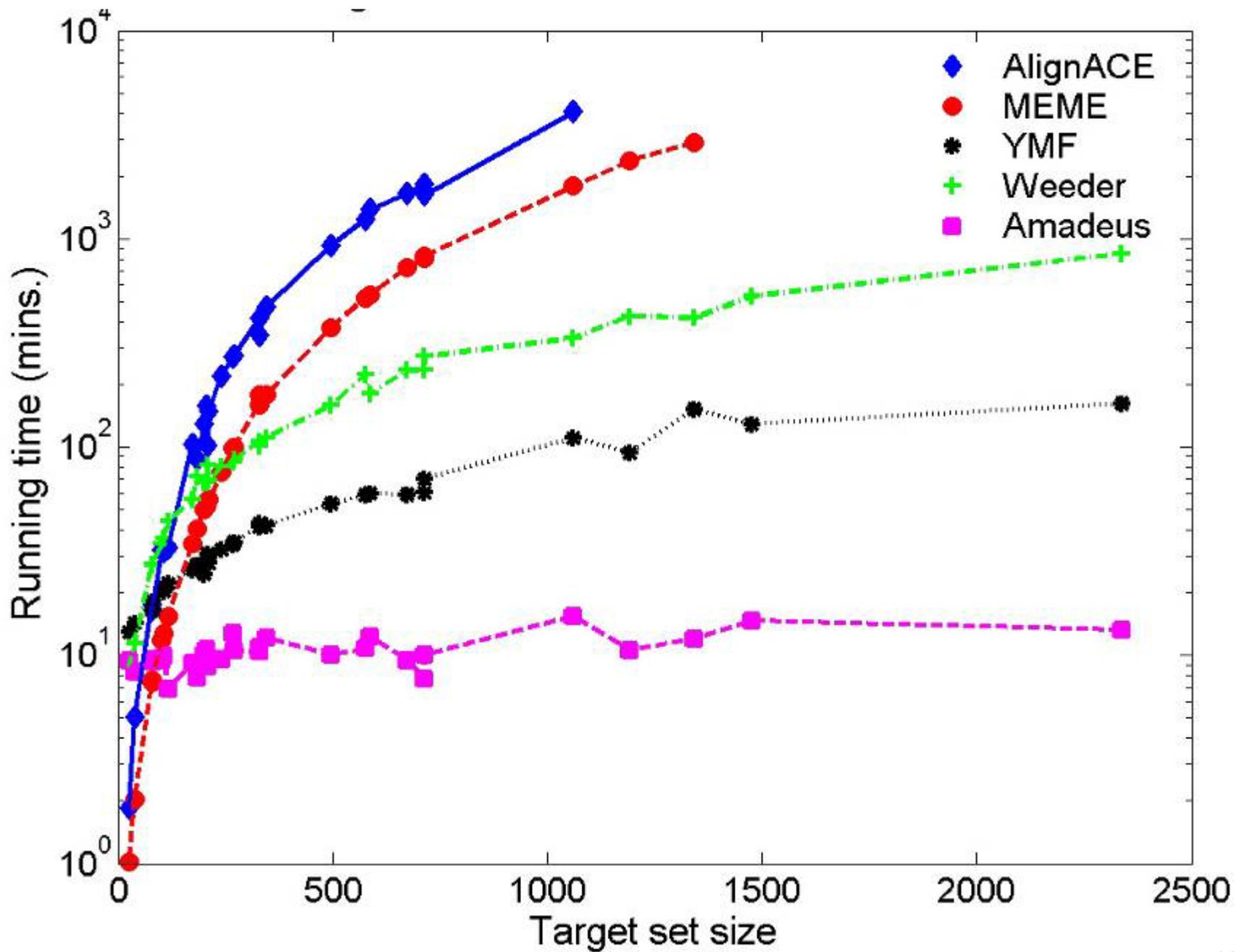
All motifs are experimentally verified

Average target set size: 400 genes (383 Kbp)



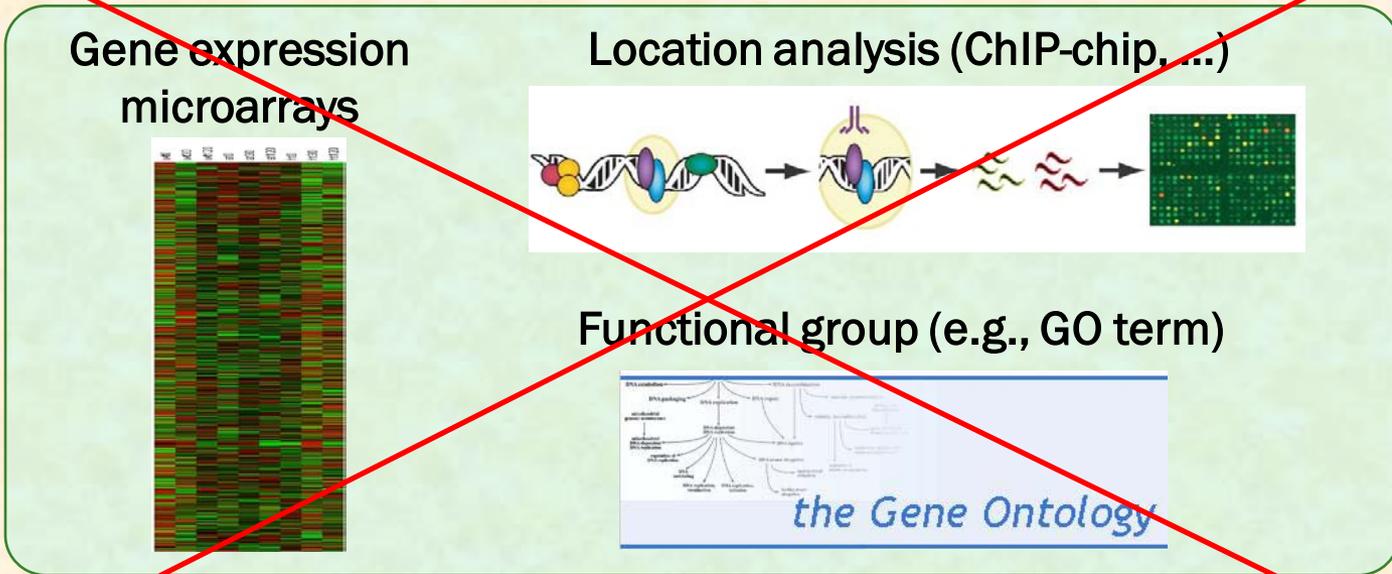
Metazoan benchmark results





Amadeus – Global spatial analysis

Co-regulated gene set



Promoter sequences



Output



Motif(s)



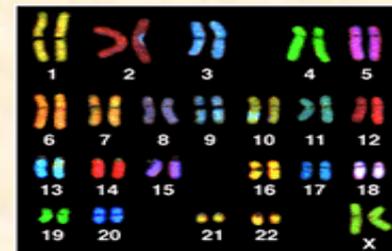
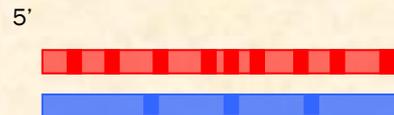
Task II: Global analyses

Scores for spatial features of motif occurrences

Input: Sequences (no target-set / expression data)

Motif scoring:

- Localization w.r.t the TSS
- Strand-bias
- Chromosomal preference





Global analysis: Chromosomal preference in *C. elegans*

Input:

- All worm promoters (~18,000)
- Score: chromosomal preference

Results:

Novel motif on chrom IV



Volume 127, Issue 6, 15 December 2006, Pages 1193-1207

Large-Scale Sequencing Reveals 21U-RNAs and Additional MicroRNAs and Endogenous siRNAs in *C. elegans*

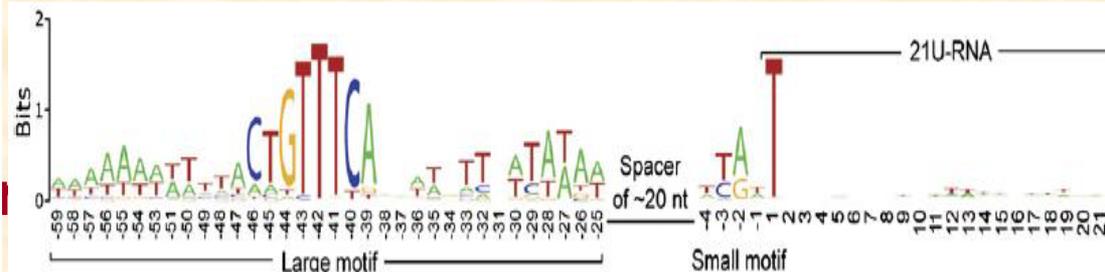
J. Graham Ruby,^{1,2} Calvin Jan,^{1,2} Christopher Player,¹ Michael J. Axtell,^{1,4} William Lee,³ Chad Nusbaum,³ Hui Ge,¹ and David P. Bartel^{1,2,*}



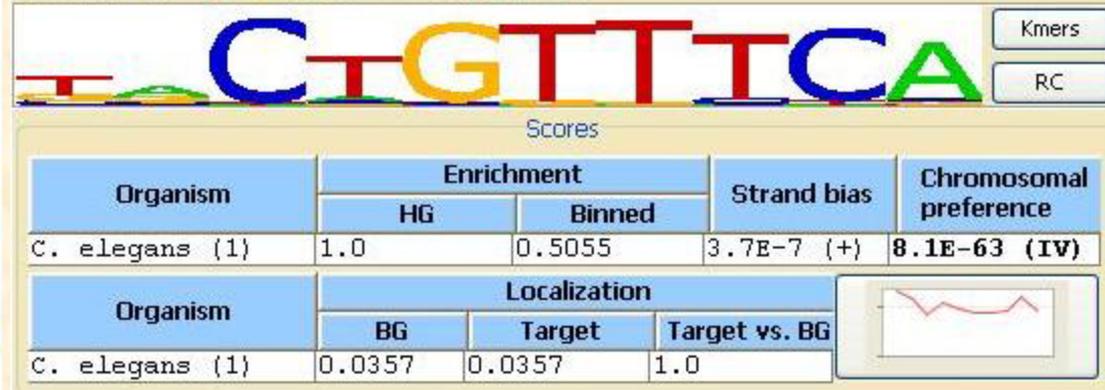
Global analysis: Chromosomal preference in *C. elegans*

SUMMARY

We sequenced ~400,000 small RNAs from *Caenorhabditis elegans*. Another 18 microRNA (miRNA) genes were identified, thereby extending to 112 our tally of confidently identified miRNA genes in *C. elegans*. Also observed were thousands of endogenous siRNAs generated by RNA-directed RNA polymerases acting preferentially on transcripts associated with spermatogenesis and transposons. In addition, a third class of nematode small RNAs, called 21U-RNAs, was discovered. 21U-RNAs are precisely 21 nucleotides long, begin with a uridine 5'-monophosphate but are diverse in their remaining 20 nucleotides, and appear modified at their 3'-terminal ribose. 21U-RNAs originate from more than 5700 genomic loci dispersed in two broad regions of chromosome IV—primarily between protein-coding genes or within their introns. These loci share a large upstream motif that enables accurate prediction of additional 21U-RNAs. The motif is conserved in other nematodes, presumably because of its importance for producing these diverse, autonomously expressed, small RNAs (dasRNAs).



Motif 1 p-value = 8.1E-63



II. Finding Transcriptional Programs



Goal

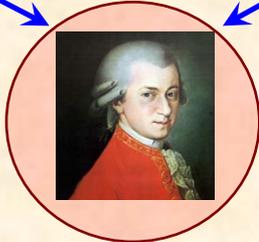
Given expression profiles, find the transcriptional programs active in them:

- the co-regulated genes,
- the motifs that govern their co-regulation



Our goal: bypass the two-step approach

Expressed data gene set



Output

1 GG 2 G 3 G 4 A 5 T 6 T 7 T 8 C 9 C 10 C

Motif(s)



Allegro: expression model

- Discretization of expression patterns

$e_1 = \text{Up (U)} \geq 1.0$

$e_2 = \text{Same (S)} (-1.0, 1.0)$

$e_3 = \text{Down (D)} \leq -1.0$

Expression pattern

	c_1	c_2	...	c_m
g	-2.3	-0.8		1.5



Discrete expression Pattern (DEP)

	c_1	c_2	...	c_m
g	D	S	...	U

- Condition frequency matrix (CFM)

$F =$

	c_1	c_2	...	c_m
U	0.05	0.1	...	0.78
S	0.9	0.2	...	0.14
D	0.05	0.7	...	0.08

- Condition weight matrix (CWM)

$$F^{(W)} = \left\{ \log \left(\frac{f_{ij}}{r_{ij}} \right) \right\} \quad (R = \{r_{ij}\} \text{ is the BG CFM})$$

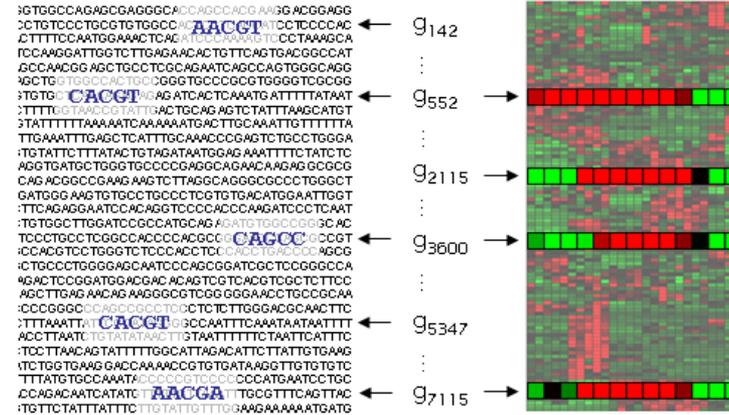
\Rightarrow Log-likelihood ratio (LLR) score

Allegro overview

Data

Cis-regulatory sequences

Gene expression matrix



Model

PWM
(sequence motif)

	p_1	p_2	p_3	...	p_k
A	1.2	6.3	-2.2		0.7
C	2.6	-0.5	4.8		0.8
G	-0.9	-1.1	0.9		-1.6
T	-0.9	-0.8	-1.9		2.6

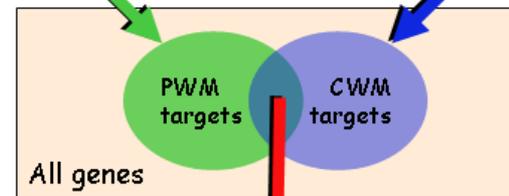
($k = 8-12$)

CWM
(expression profile)

	c_1	c_2	c_3	...	c_m
Up	0.2	0.8	4.2		-4.5
Same	-3.2	-1.1	0.3		-0.3
Down	9.5	3.2	-1.2		9.1

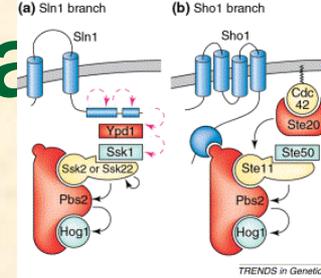
($m = 1-150+$)

Model evaluation & optimization

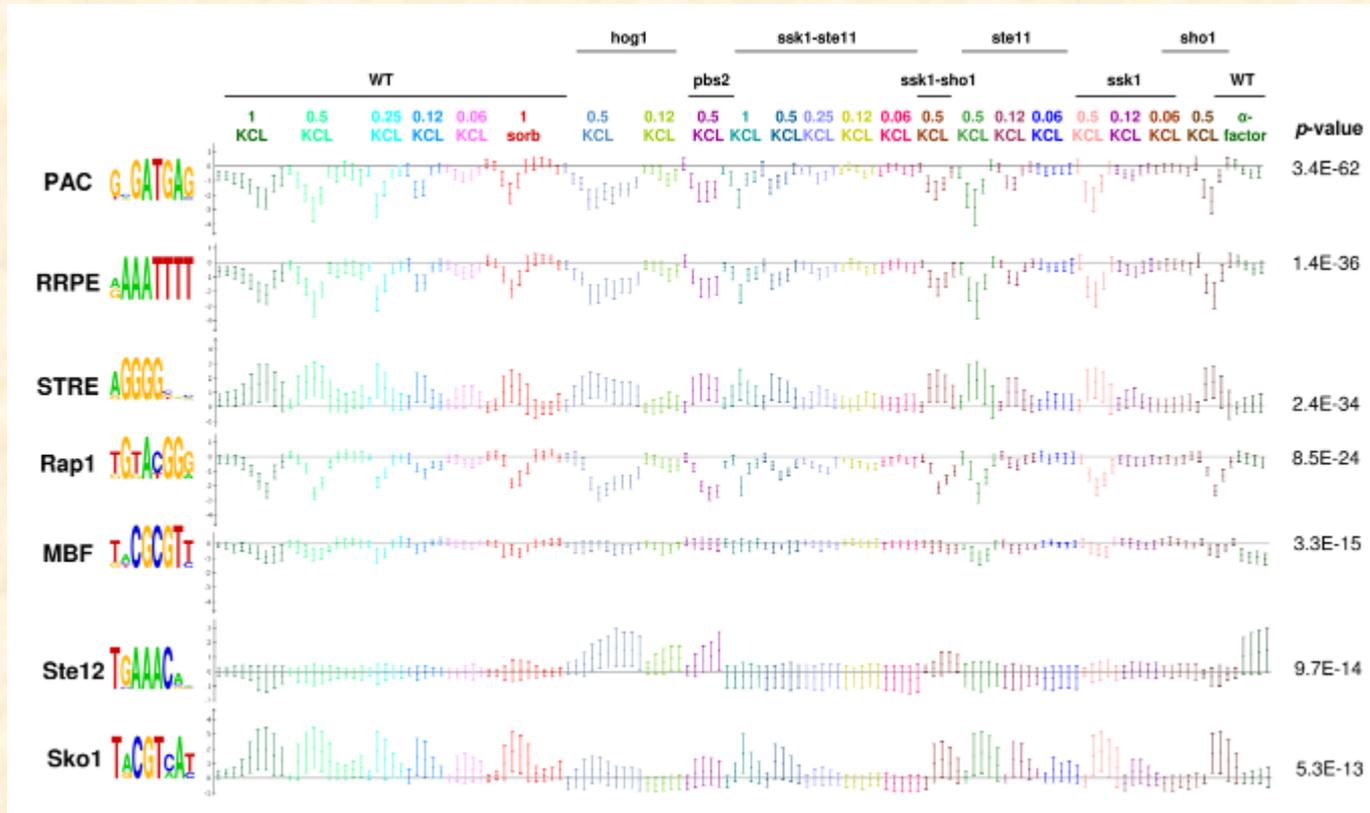


Score = Enrichment p -value

Yeast osmotic shock pathway



- ~6,000 genes, 133 conditions [O'Rourke et al. '04]



- Allegro can discover **multiple motifs** with **diverse expression patterns**, even if the response is in a **small fraction of the conditions**
- Extant two-step techniques** recovered only 4 of the above motifs:
 - K-means/CLICK + Amadeus/Weeder: RRPE, PAC, MBF, STRE
 - Iclust + FIRE: RRPE, PAC, Rap1, STRE



3' UTR analysis: Human stem cells

- ~14,000 genes, 124 conditions (various types of proliferating cells) [Mueller et. al, Nature'08]
- Biases in length / GC-content of 3' UTRs, e.g.:

<u>100 highly-expressed genes in...</u>	3' UTR: <u>length</u>	<u>GC</u>
Embryoid bodies	584	47%
Undifferentiated ESCs	774	44%
ESC-derived fibroblasts	1240	39%
Fetal NSCs	1422	43%

(ESCs = embryonic stem cells, NSCs = neural stem cells)

- Extant methods / Allegro with HG score: report only false positives



Human stem cells: results using binned score

Current knowledge

Most highly expressed miRNAs in human/mouse ESCs

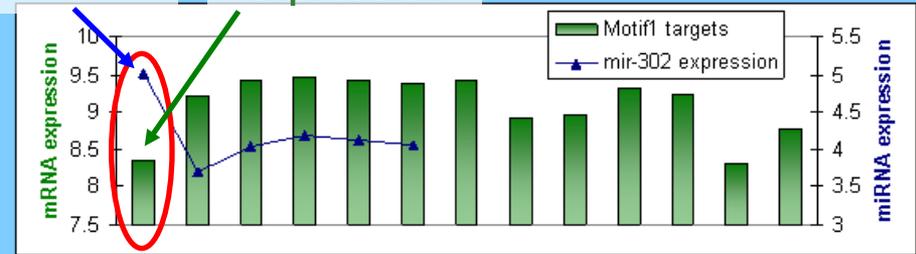
Motif 1: $p\text{-value} = 2 \cdot 10^{-13}$



hsa-miR-302a	UAAGUGCUUCC...
hsa-miR-17	CAAAGUGCUUA...
hsa-miR-372	AAAGUGCUGCG...
hsa-miR-520e	AAAGUGCUUCC...

miRNA expression

targets expression

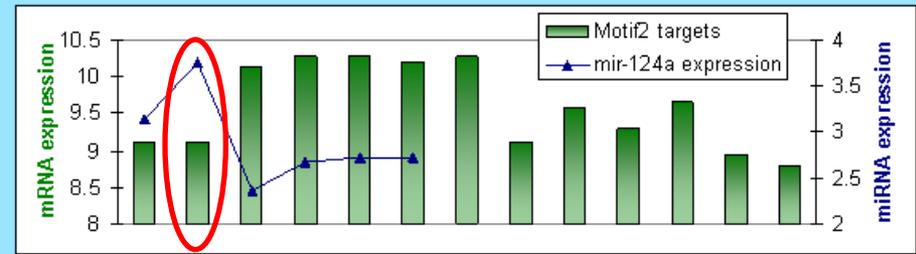


Motif 2: $p\text{-value} = 7 \cdot 10^{-13}$



hsa-miR-124	UAAGGCACGCG...
hsa-miR-506	UAAGGCACCCU...

Abundant & functional in neural cell lineage

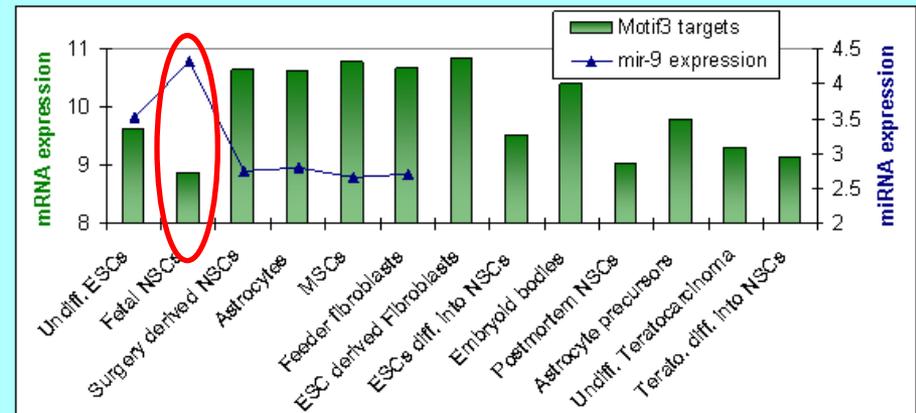


Motif 3: $p\text{-value} = 10^{-9}$



hsa-miR-9	UCUUUGGUUAU...
-----------	----------------

Expressed specifically in neural lineage; active role in neurogenesis





Chaim Linhart



Yonit Halperin



Igor Ulitsky



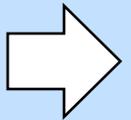
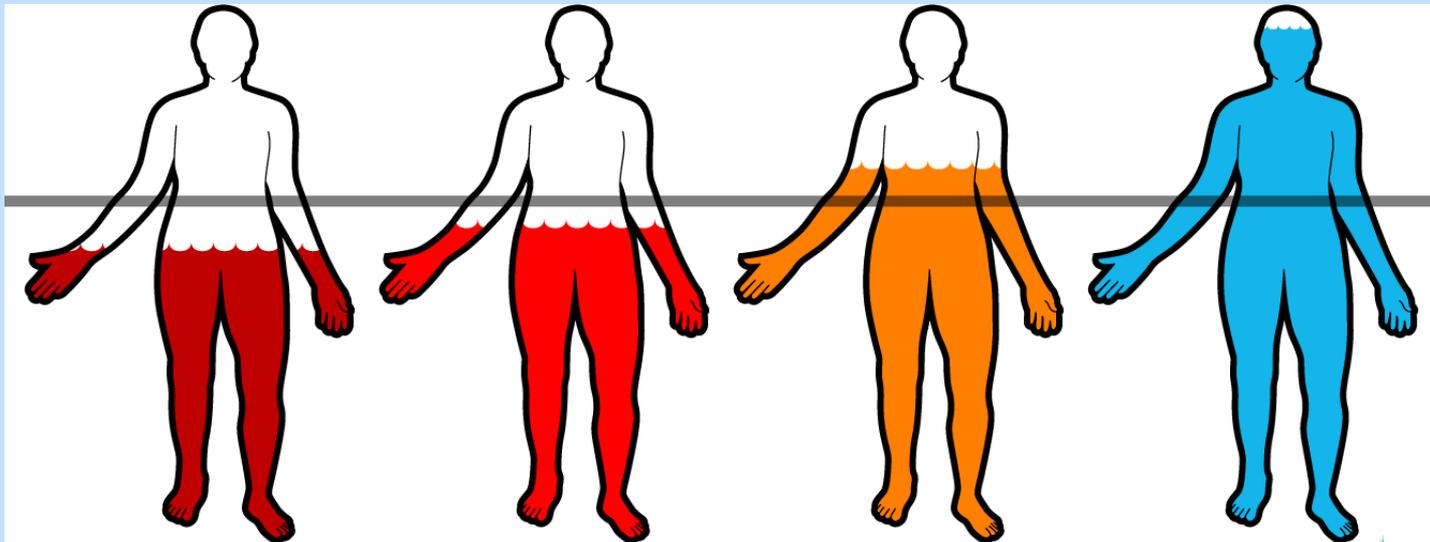
Yaron
Orenstein

Open questions

- Better PWM inference: new scores, algs
- Richer models for in vivo / in vitro data – really helpful or diminishing return?
- How to evaluate model quality: match to literature? Ranking based? In vivo? In vitro?
- Integration of motif finding & expression
- Principled means to find motif pairs



Using expression profiles and protein networks to understand cancer I

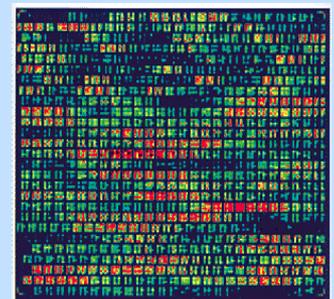
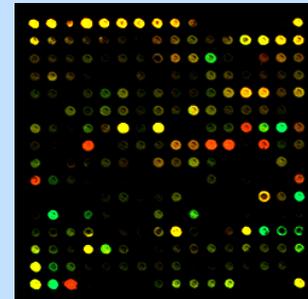


DNA chips / Microarrays

- Simultaneous measurement of expression levels of all genes.
- Global view of cellular processes.
- > 800,000 profiles available in ArrayExpress

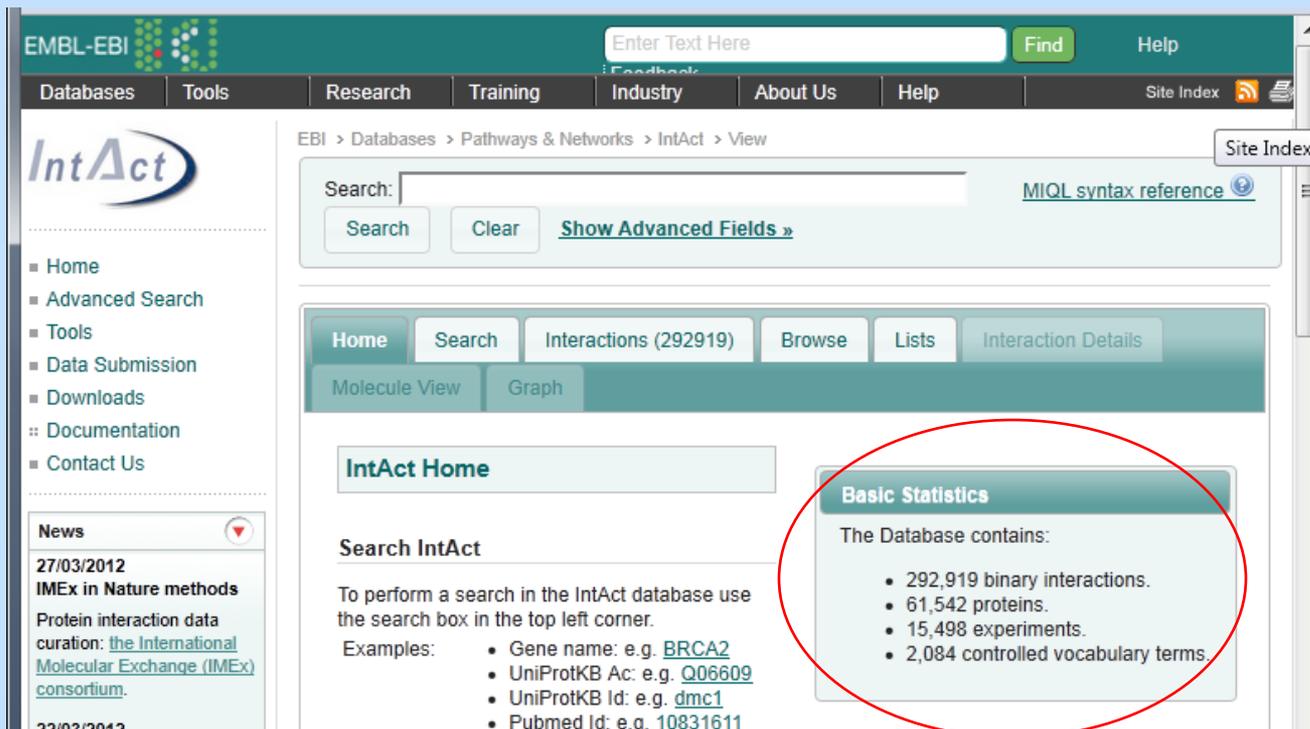


Affymetrix



Protein-protein interactions (PPIs)

- A regulates/binds to B
- High throughput: abundant, noisy
- Large, readily available resource

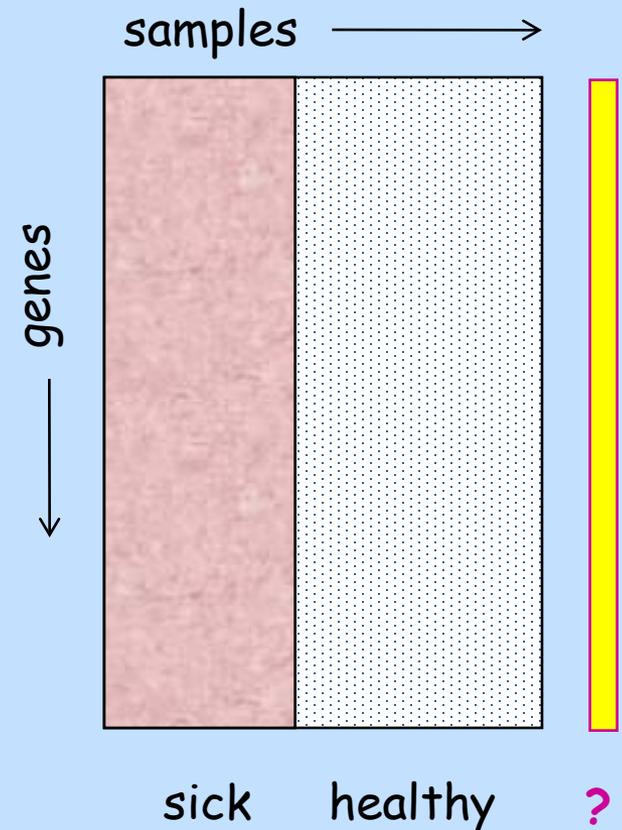


The screenshot shows the IntAct database interface. At the top, there is a search bar with the text "Enter Text Here" and a "Find" button. Below the search bar is a navigation menu with options: Databases, Tools, Research, Training, Industry, About Us, Help, Site Index, and RSS. The main content area features the IntAct logo and a breadcrumb trail: "EBI > Databases > Pathways & Networks > IntAct > View". A search box is present with "Search" and "Clear" buttons, and a link to "Show Advanced Fields »". Below this is a secondary navigation bar with tabs for "Home", "Search", "Interactions (292919)", "Browse", "Lists", and "Interaction Details". Underneath, there are options for "Molecule View" and "Graph". The "Basic Statistics" section, highlighted with a red circle, states: "The Database contains:" followed by a list:

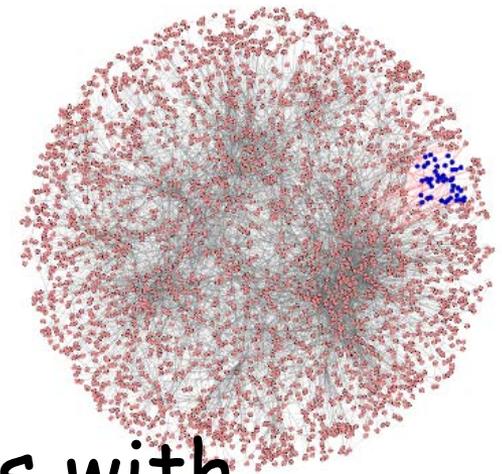
- 292,919 binary interactions.
- 61,542 proteins.
- 15,498 experiments.
- 2,084 controlled vocabulary terms.

Case/control studies

- A typical study: 100s expression profiles of sick (**case**) & healthy (**control**) individuals
- **Classification:** Given a partition of the samples into types, classify the types of new samples
- Can the network help?



The network angle

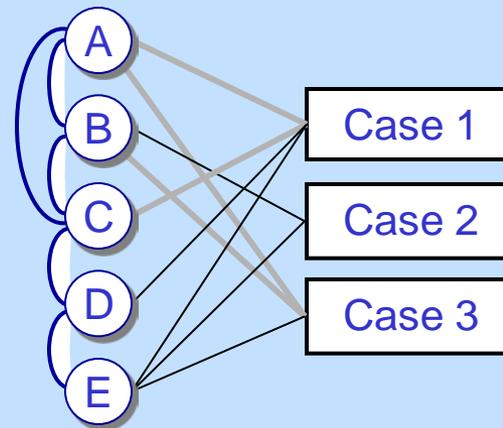


- Integrate case-control profiles with network information
- Extract dysregulated pathways **specific** to the cases
- Account for **heterogeneity** among cases
- Meaningful pathway: **connected**



Preprocessing

- For each gene, use the distribution of values among the controls to decide if the gene is **dysregulated** in each of the cases



	Control 1	Control 2	Control 3	Control 4	Case 1	Case 2	Case 3
A	Yellow	Light Yellow	Yellow	Pink	Dark Blue	Light Blue	Dark Blue
B	Yellow	Light Yellow	Yellow	Yellow	Light Yellow	Dark Blue	Dark Blue
C	Light Green	Yellow	Light Green	Light Green	Dark Blue	Light Grey	Light Purple
D	Light Green	Yellow	Yellow	Yellow	Dark Blue	Purple	Light Green
E	Yellow	Light Yellow	Yellow	Light Yellow	Dark Blue	Dark Blue	Dark Blue

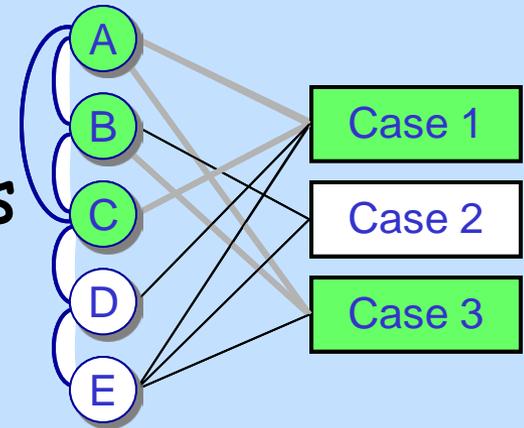
↓

	Case 1	Case 2	Case 3
A	1	0	1
B	0	1	1
C	1	0	0
D	1	0	0
E	1	1	1

←

Dysregulated pathway

- Input:
 - Bipartite graph: genes, cases
 - Edge (gene g , case c) if g is dysregulated in c
 - A network over the genes
- **Dysregulated pathway (DP):** smallest connected subnetwork s.t. sufficiently many $\geq k$ genes are dysregulated in all but few $\leq l$ cases
- Small pathway \rightarrow focused disease explanation



$$k=2, l=1$$



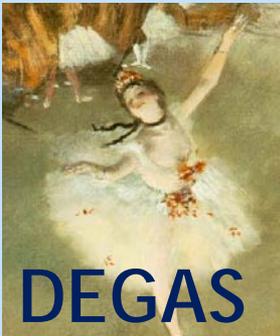
Min connected set cover problem

Complexity

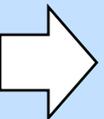
- **Set cover** problem: Given sets of elements, find fewest sets that cover all elements

k	l	G	Problem
1	0	Clique	Set cover
k	0	Clique	Set k-cover
1	>0	Clique	Partial set cover
1	0	Any	Connected set cover (Shuai & Hu 06)

- All are NP-Hard
- Devised approximation and heuristic algs



DysrEgulated Gene set
Analysis via Subnetworks



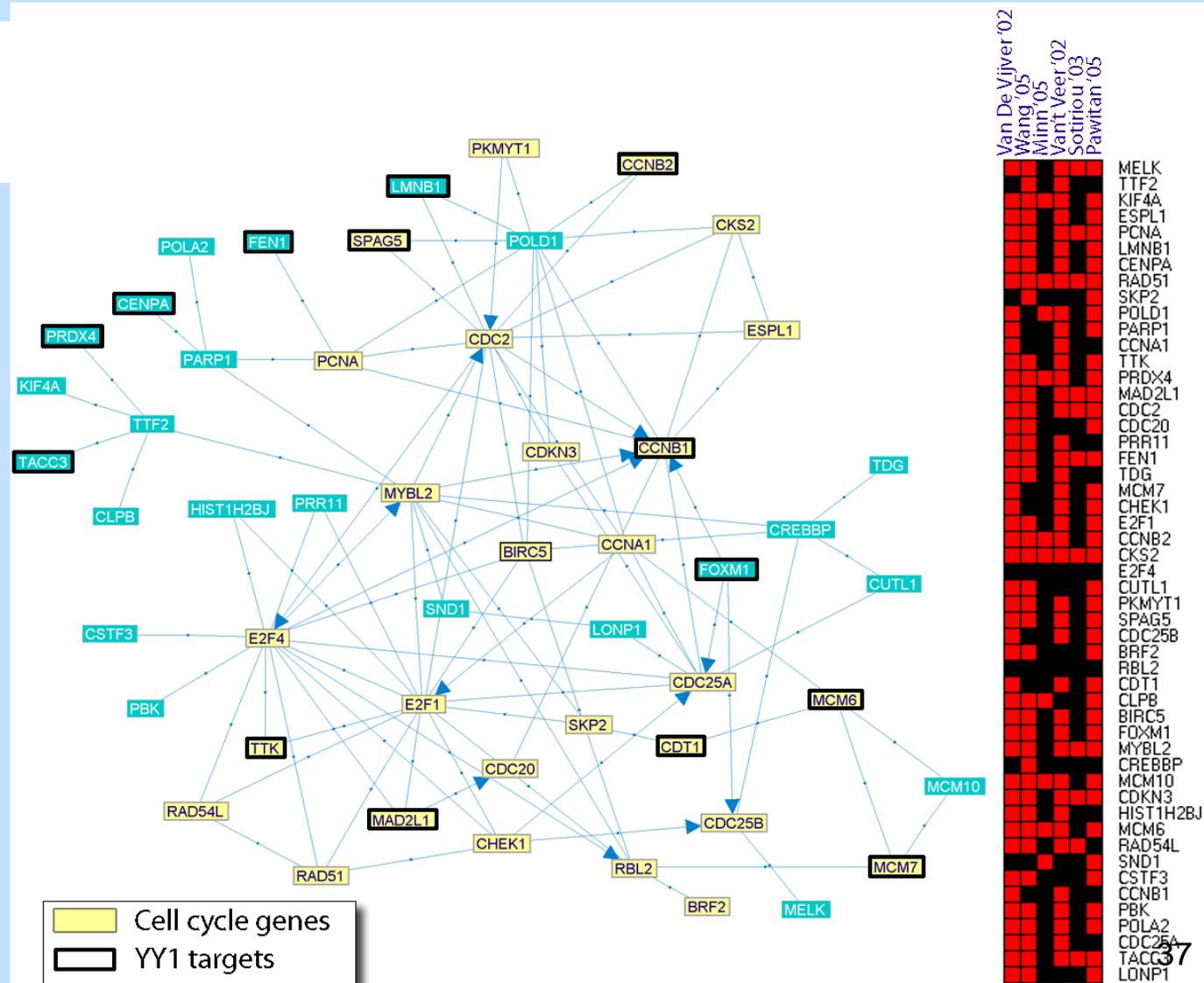
Breast cancer meta-analysis

- 6 breast cancer studies comparing poor and good prognosis
 - Van't Veer *et al. Nature* 2002
 - Van de Vijver *et al. NEJM* 2002
 - Wang *et al. Lancet* 2005
 - Minn *et al. Nature* 2005
 - Sotiriou *et al. PNAS* 2003
 - Pawitan *et al. Breast Cancer Research* 2005
- Poor prognosis = metastases within 5 years
- 1,004 patients in total
- Elements = studies
- Discovered 2 significant DPs associated with **poor** prognosis and one associated with **good** prognosis



Poor prognosis network 1

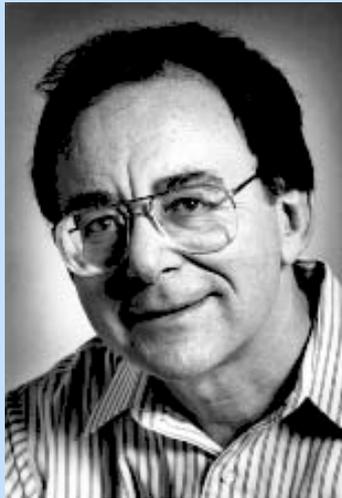
- $k = 40, l = 2,$
 $p < 0.005$
- Enriched with cell-cycle associated genes ($p = 2 \cdot 10^{-26}$) & YY1 targets ($p = 2.42 \cdot 10^{-16}$)
- Enriched with genes localized to the nucleus



Summary

- A method for finding subnetworks of dysregulated genes
- Specific to cases, but allows outliers and exception
- Connected set cover paradigm
- Better approximations??





Dick Karp,
Berkeley



Igor Ulitsky,
Whitehead Inst



Akshay
Krishnamurthy
CMU

T H A N K Y O U
1 2 3 4 5 6 7 8



postdocs available

Support: Israel Academy of Sciences, Wolfson Foundation, Edmond J. Safra Foundation, US-Israel BSF, German-Israeli Fund, EU 6th and 7th Frameworks, Intel, IBM, I-CORE Gene regulation & disease.