

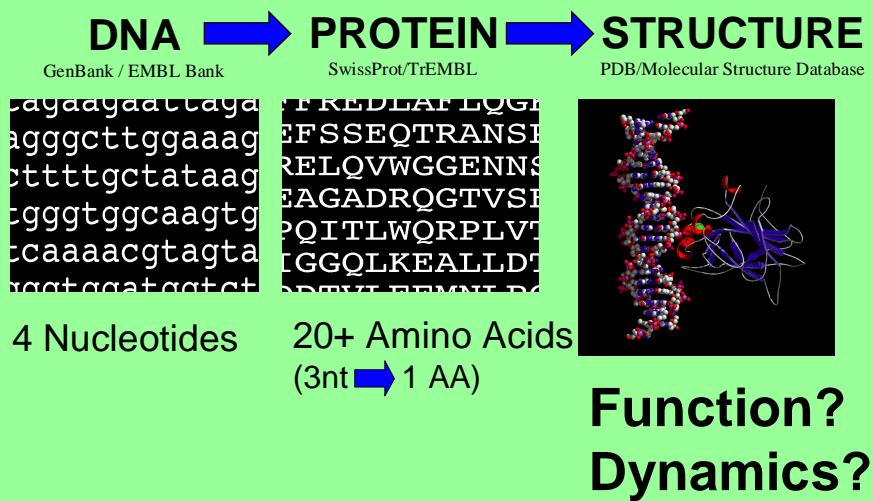
# Studying gene regulation by data mining approaches

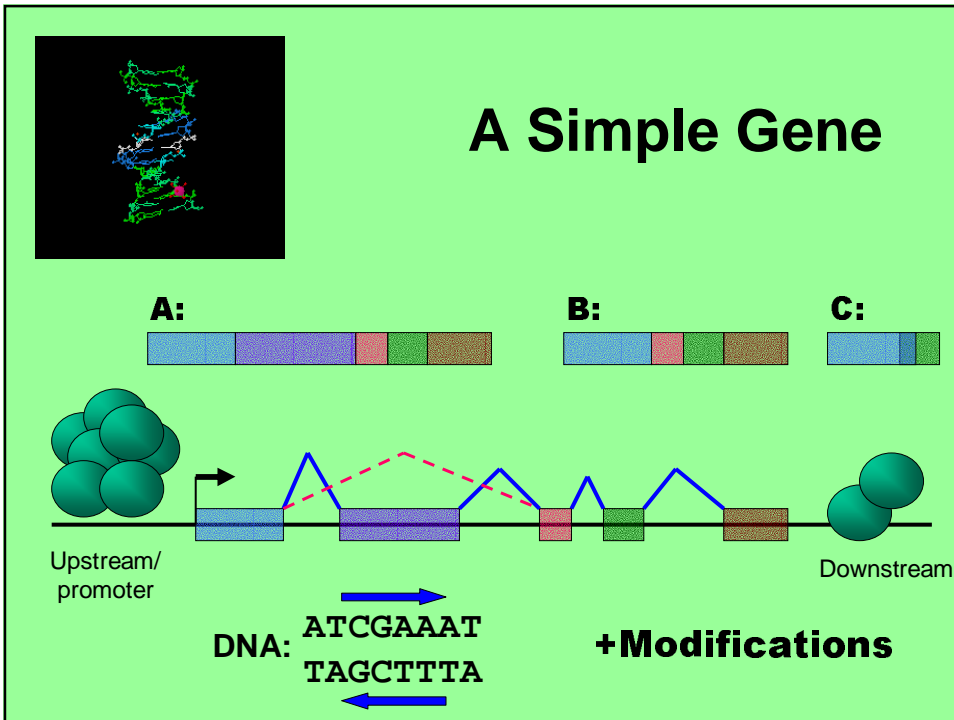
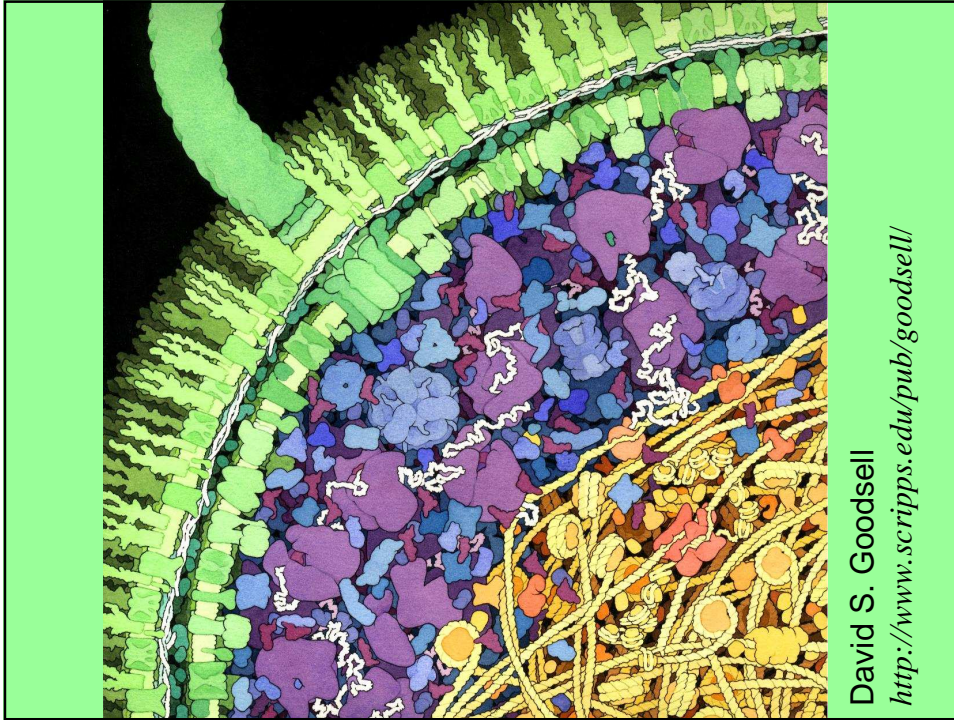
Jaak Vilo  
vilo@egeen.ee



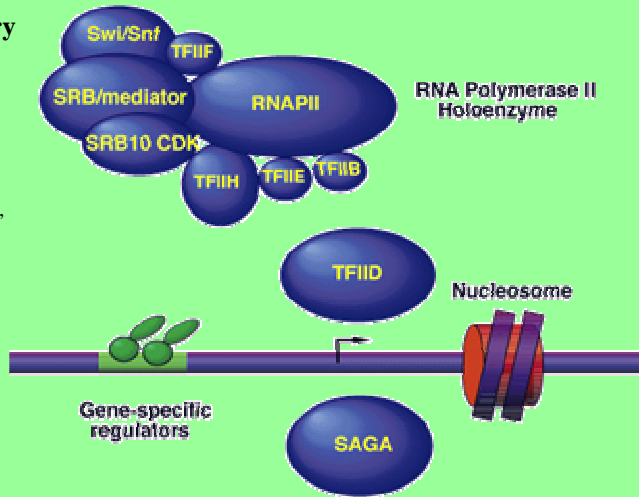
Estonian Computer Science Theory days: Pedase, 3.10.2003

## DNA determines function?





Model of RNA Polymerase II Transcription Initiation Machinery. **The machinery depicted here encompasses over 85 polypeptides in ten (sub) complexes:** core RNA polymerase II (RNAPII) consists of 12 subunits; TFIID, 9 subunits; TFIIE, 2 subunits; TFIIH, 3 subunits; TFIIB, 1 subunit, TFIID, 14 subunits; core SRB/mediator, more than 16 subunits; Swi/Snf complex, 11 subunits; Srb10 kinase complex, 4 subunits; and SAGA, 13 subunits.



F.C.P. Holstege, E.G. Jennings, J.J. Wyrick, Tong Ihn Lee, C.J. Hengartner, M.R. Green, T.R. Golub, E.S. Lander, and R.A. Young  
Dissecting the Regulatory Circuitry of a Eukaryotic Genome  
Cell 95: 717-728 (1998)

```

TGTTCTTTCTTCTTTTCATACATCCTTTTCCTTTTTTTCC
TTCTCCTTTCATTTCTGACTTTTAATATAGGCTTACCA
TCCTTCTTCTCTTCAATAACCTTCTTACATTGCTTCTTC
TTCGATTGCTTCAAAGTAGTTCGTGAATCATCCTTCAAT
GCCTCAGCACCTTCAGCACTTGCACTTCATTCTCTGGAA
GTGCTGCACCTGCGCTGTCTTGCTAATGGATTTGGAGTT
GGCGTGGCACTGATTTCTTCGACATGGCGGCGTCTTCT
TCGAATTCCATCAGTCCCTCATAGTTCTGTTGGTTCTTTT
CTCTGATGATCGTCATCTTTCACTGATCTGATGTTCTTG
TGCCCTATCTATATCATCTCAAAGTTCACCTTTGCCACT
TTCCAAGATCTCTCATTCAATAATGGGCTTAAAGCCGTAC
TTTTTTCACCTCGATGAGCTATAAGAGTTTTTCACTTTTA
GATCGTGGCTGGGCTTATATTACGGTGTGATGAGGGCGC
TTGAAAAGATTTTTTTCATCTCACAAGCGACGAGGGCCCG
AGTGTGTTGAAGCTAGATGCAGTAGGTGCAAGCGTAGAGT
CTTAGAAGATAAAGTAGTGAATTACAATAGATTCGATAC

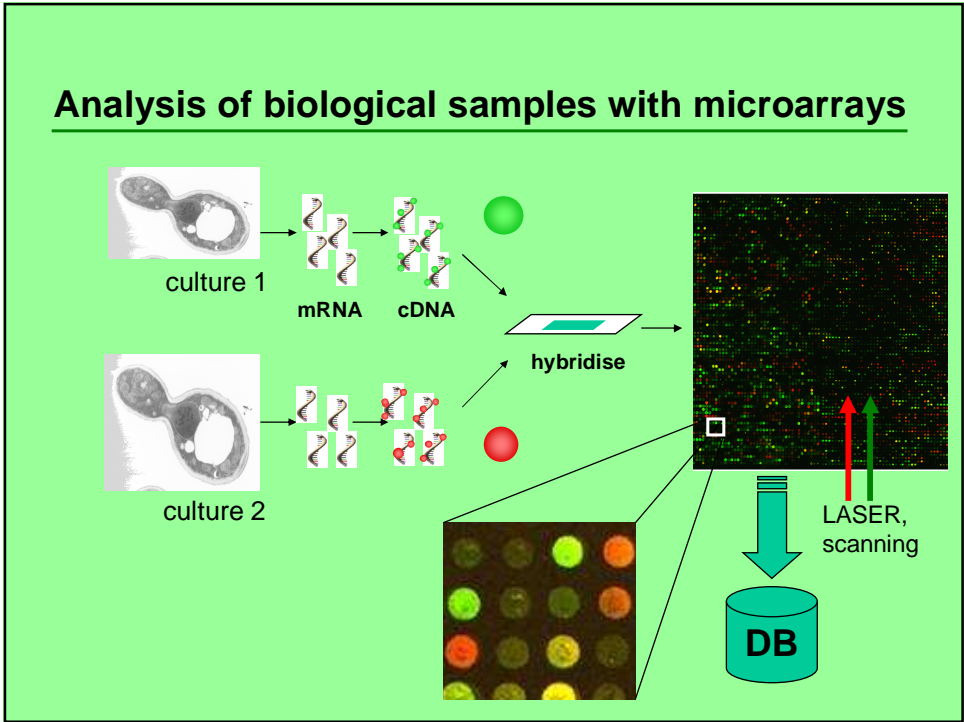
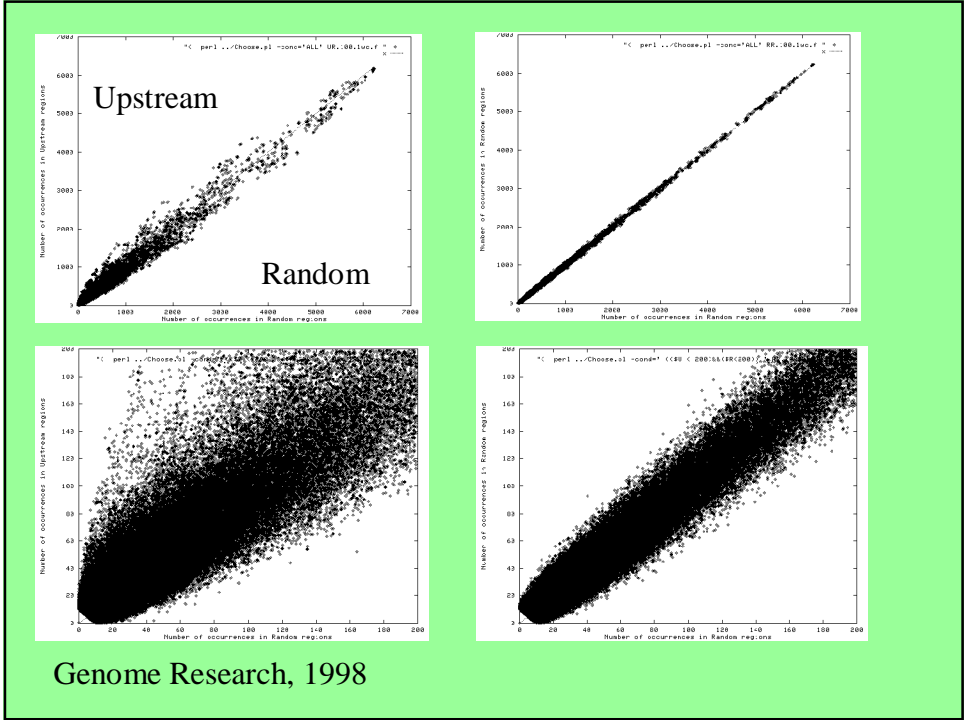
```

## Patterns: AT

```
TGTTCTTTCTTCTTTATACATCCTTTTCCTTTTTTTCC
TTCTCCTTTATTTCCCTGACTTTTAAATATAGGCTTACCA
TCCTTCTTCTCTTCAATAACCTTCTTACATTGCTTCTTC
TTCGATTGCTTCAAAGTAGTTCGTGAATCATCCTTCAAT
GCCTCAGCACCTTCAGCACTTGCACTTCATTCTCTGGAA
GTGCTGCACCTGCGCTGTCTTGCTAATGGATTTGGAGTT
GGCGTGGCACTGATTTCTTCGACATGGCGGGCGTCTTCT
TCGATTCCATCAGTCCATAGTTCTGTTGGTTCTTTT
CTCTGATGATCGTCATCTTTCACTGATCTGATGTTCCCTG
TGCCCTATCTATATCATCTCAAAGTTCACTTTGCCACT
TTCCAAGATCTCTATTCATATGGGCTTAAGCCGTAC
TTTTTCACTCGATGAGCTATAAGAGTTTTCCACTTTTA
GATCGTGGCTGGGCTTATATTACGGTGTGATGAGGGCGC
TTGAAAAGATTTTTTCACTCACAAGCGACGAGGGCCCCG
AGTGTTTGAGCTAGATGCAGTAGGTGCAAGCGTAGAGT
CTTAGAAGATAAAGTAGTGAATTACAATAGATTCGATAC
```

## Patterns: [AT][ACT]AT (WHAT)

```
TGTTCTTTCTTCTTTCATACATCCTTTTCCTTTTTTTCC
TTCTCCTTTCATTTCCCTGACTTTTAATATAGGCTTACCA
TCCTTCTTCTCTTCAATAACCTTCTTACATTGCTTCTTC
TTCGATTGCTTCAAAGTAGTTCGTGAATCATCCTTCAAT
GCCTCAGCACCTTCAGCACTTGCACTTCATTCTCTGGAA
GTGCTGCACCTGCGCTGTCTTGCTAATGGATTTGGAGTT
GGCGTGGCACTGATTTCTTCGACATGGCGGGCGTCTTCT
TCGATTCCATCAGTCCTCATAGTTCTGTTGGTTCTTTT
CTCTGATGATCGTCATCTTTCACTGATCTGATGTTCCCTG
TGCCCTATCTATATCATCTCAAAGTTCACTTTGCCACT
TTCCAAGATCTCTTCATTCATATGGGCTTAAGCCGTAC
TTTTTCACTCGATGAGCTATAAAGAGTTTTCCACTTTTA
GATCGTGGCTGGGCTTATATTACGGTGTGATGAGGGCGC
TTGAAAAGATTTTTTTCATCTCACAAGCGACGAGGGCCCCG
AGTGTTTGAGCTAGATGCAGTAGGTGCAAGCGTAGAGT
CTTAGAAGATAAGTAGTGAATTACAATAGATTCGATAC
```



## From microarray images to gene expression data

Raw data

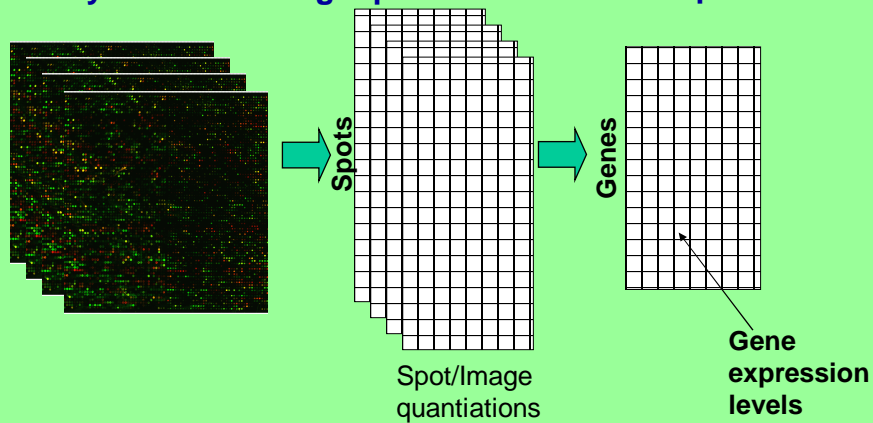
Intermediate data

Final data

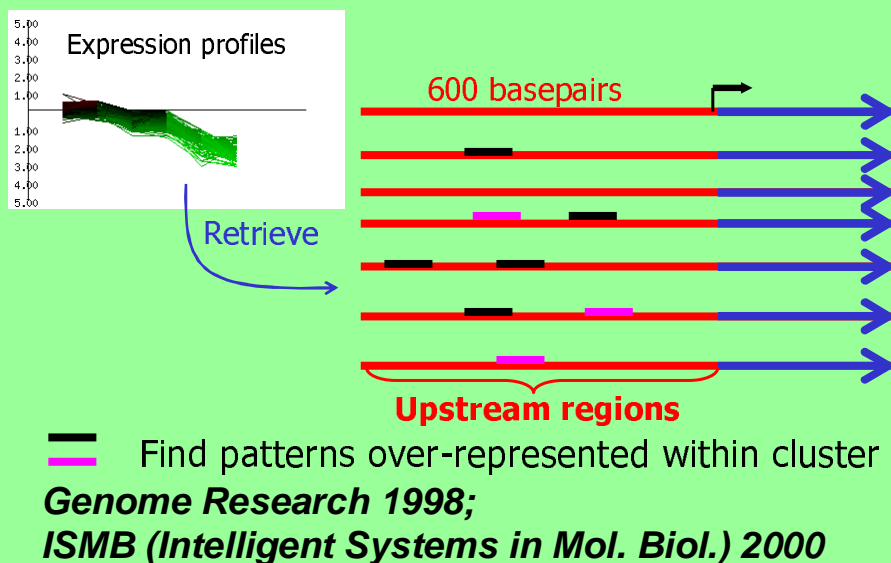
Array scans

Image quantifications

Samples



## Cluster of co-expressed genes, pattern discovery in regulatory regions

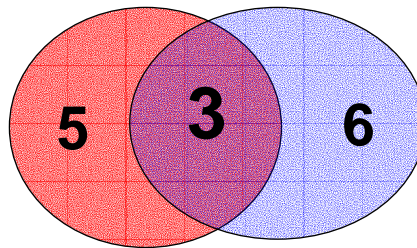




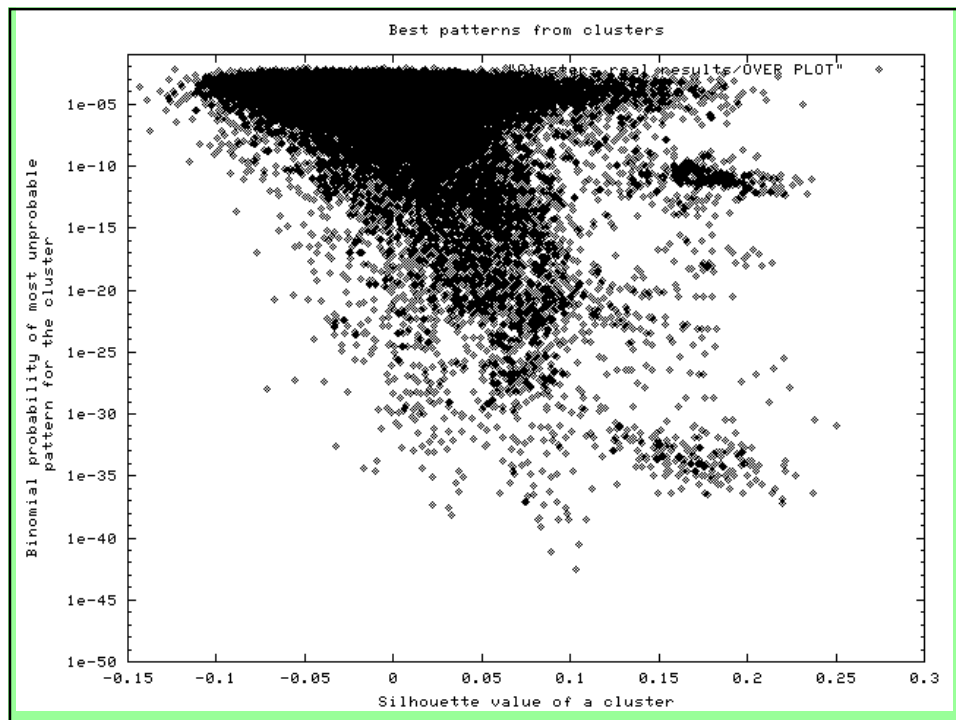


## Set overlap

25 genes

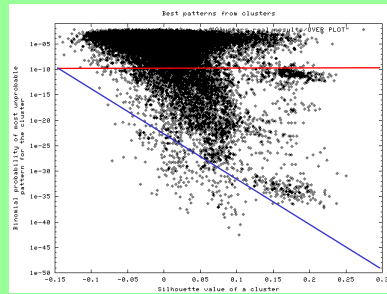


$P(\text{choose 6 balls randomly from 25, of which 5 reds, and observe 3 or more red})$

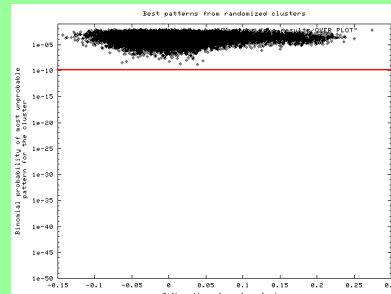




## Pattern vs cluster “strength”



The pattern probability vs. the average silhouette for the cluster



The same for randomised clusters

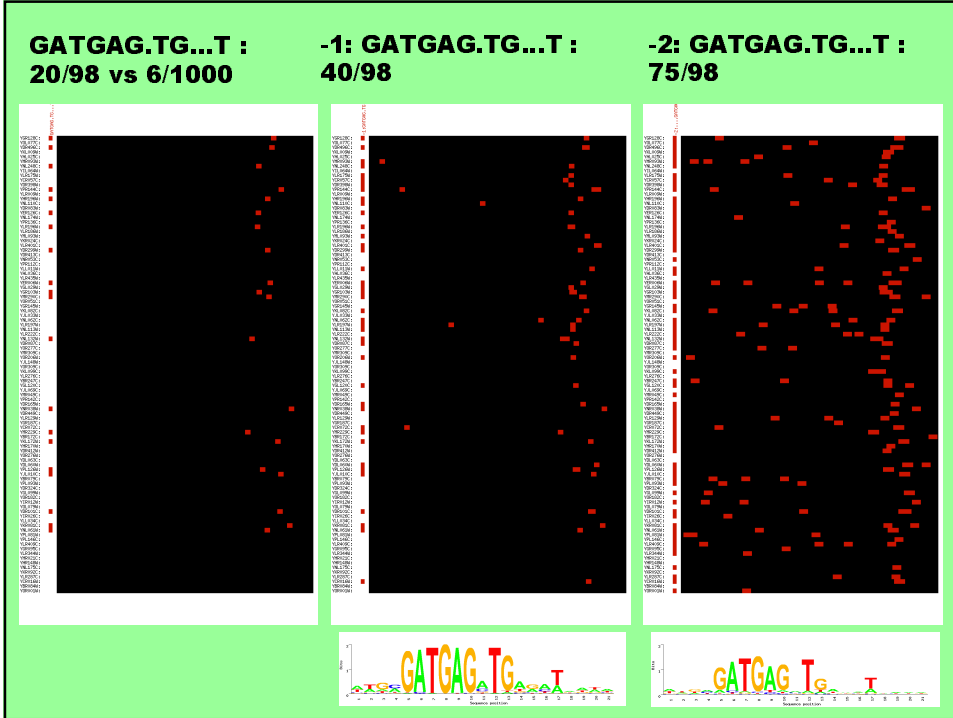
Vilo *et.al.* ISMB 2000

## Regular patterns (SPEXS)

- Substrings ATCGA
- Add groups ATC[GC][AT]
- Add (unrestricted) wildcards AT\*CG
- Add restricted wildcards AT\*(2,5)CG
- Combine all above

**AT[GC]\*(1,3)[GT]AC**

**TGC.....ACG**



## Consensus matrix building

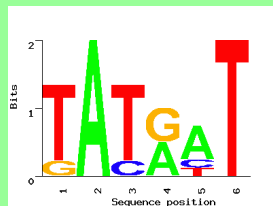
TACGAT  
TATAAT  
TATAAT  
GATACT  
TATGAT  
TATGTT

A	0	6	0	3	4	0
C	0	0	1	0	1	0
G	1	0	0	3	0	0
T	5	0	5	0	1	6

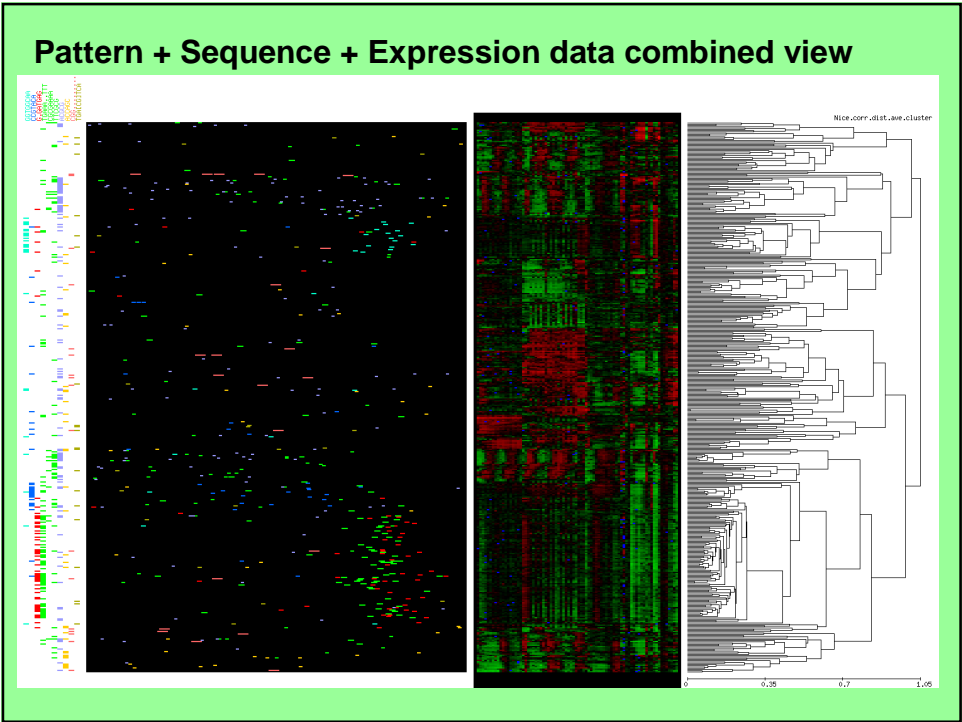
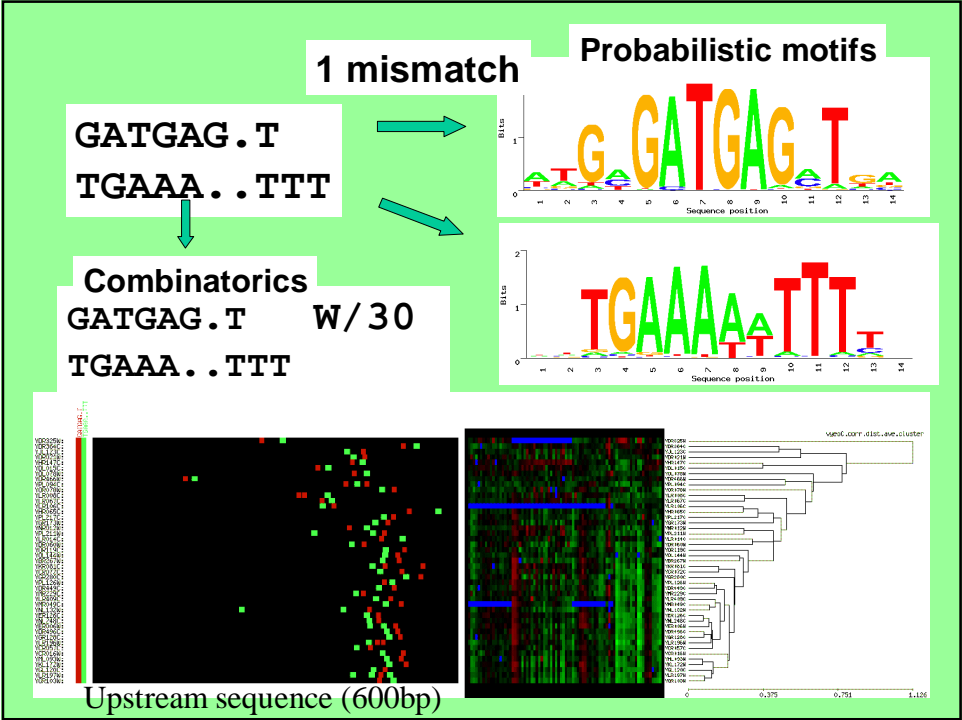
$$I_i = 2 + \sum_{b \in A,C,T,G} f_{b,i} \log_2 f_{b,i}$$

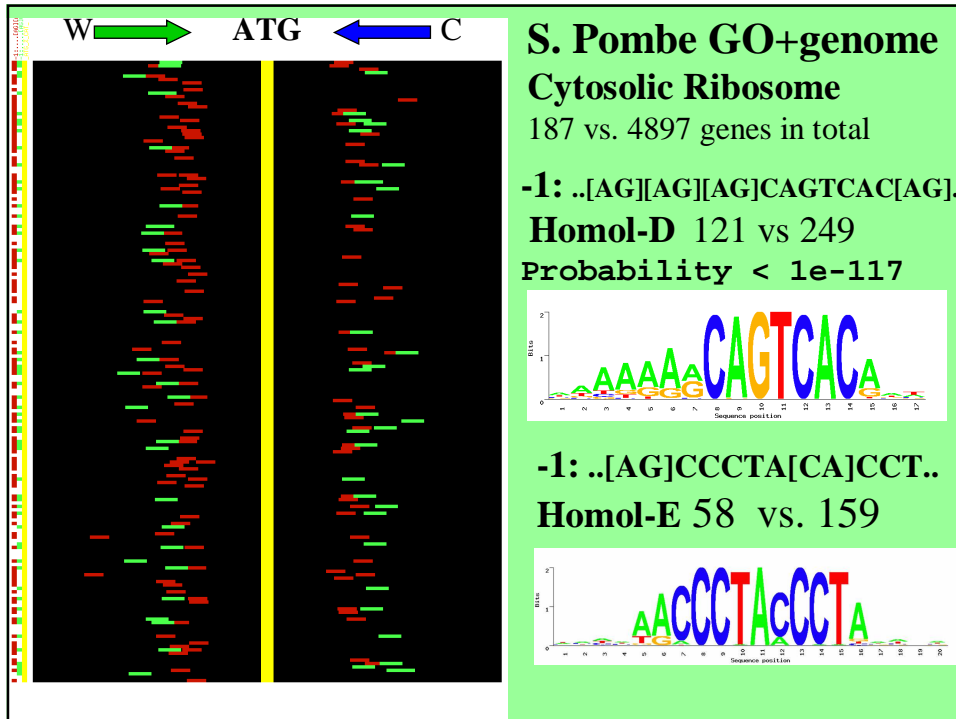
Consensi:

TATAAT  
TATRNT  
[GT]A[CT][AG][ACT]T



$$-f_{b,i} \log_2 f_{b,i}$$





## SPEXS - Sequence Pattern EXhaustive Search

Jaak Vilo, 1998, 2002

- **User-definable pattern language:** substrings, character groups, wildcards, flexible wildcards (c.f. **PROSITE**)
- Fast exhaustive search over pattern language
- “Lazy suffix tree construction”-like algorithm
- **Analyze multiple sets of sequences simultaneously**
- Restrict search to most frequent patterns only (in each set)
- **Report** most frequent patterns, patterns over- or underrepresented in selected subsets, or patterns significant by various statistical criteria, e.g. by binomial distribution

## Suffix tree – represent all suffixes

CATAT => suffix tree

123456

CATAT\$ 1

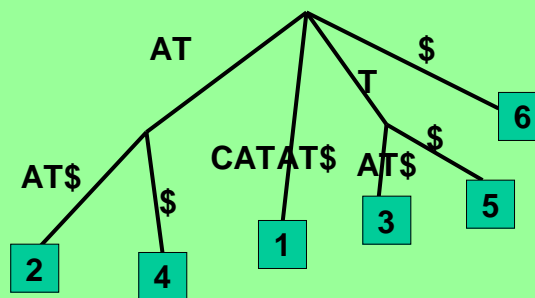
ATAT\$ 2

TAT\$ 3

AT\$ 4

T\$ 5

\$ 6



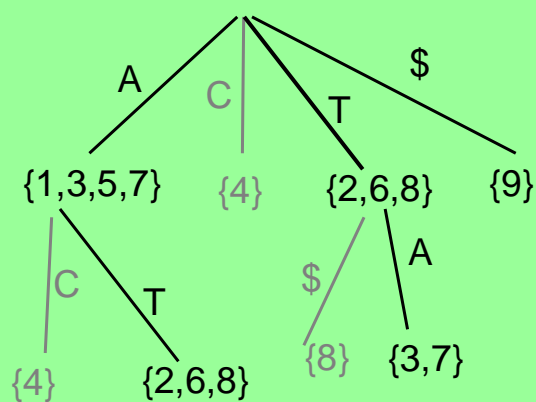
O(n) time and space

## “Lazy” construction of trie

ATACATAT\$

123456789  
ATACATAT\$

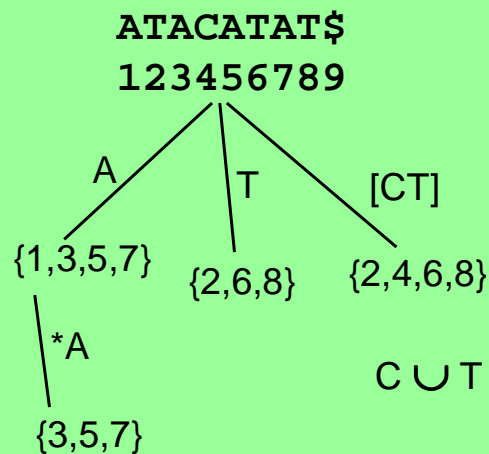
- Suffix trie
- O(n<sup>2</sup>)
- Kurtz, Giegerich
- Good in practice



....

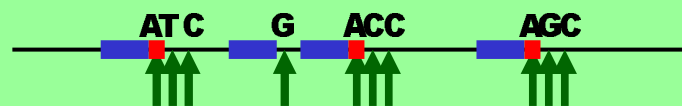
## SPEXS: pattern discovery based on pattern trie.

- Substrings
- Group characters
- Wildcard positions
- Variable length wildcards
- Restrictions on the number on each separately
- **At least k occurrences**
- Exact occurrences locations for each pattern



Vilo 1998, 2002

## Sequence patterns: the basis of the SPEXS



- GCAT (4 positions)
- GCATA (3 positions)
- GCATA.
- GCATA.C

## SPEXS: specify the pattern language and parameters for pattern discovery

**Sequences**

**Background**

**Pattern language**

**“Fitness”**

**Search order**

**Pattern frequency**

## Combinatorics of sites

- Which binding sites tend to co-occur frequently together in upstreams
- Association rules data mining
  - $\#(A,B) = 200$  ,  $\#(A,B,C) = 180$
  - $A,B \Rightarrow C$  (90%)

- *Alvis Brazma, Jaak Vilo, Esko Ukkonen and Kimmo Valtonen*  
[Data Mining for Regulatory Elements in Yeast Genome.](#)  
 Fifth International Conference on Intelligent Systems for Molecular Biology, ISMB-97 (pp. 65-74) June, 1997. AAAI Press.



## Research goals

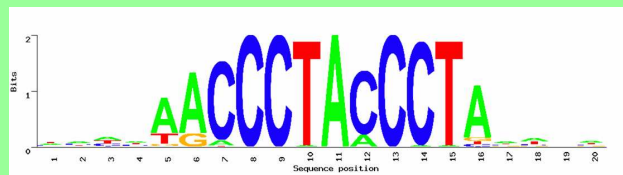
- **Generate a (full) list of hypothetical regulatory signals (for each and every gene)**
- **Maintain a DB of all known or predicted inf.**
- How does this correlate to known information and/or experimental data (e.g. ChIP on chip)
- Predict from unannotated DNA where are the promoters (and genes)
- Predict from DNA how the gene is expressed given concentrations of all TF-s in cell
- Predict the alternative splicing isoforms
- Evolution & comparative genomics approaches

## How to know what is known?

- After in silico predictions the first question should be
  - **How does that compare to current knowledge?**
- But what if databases do not allow to answer such questions easily?

Similarity?  
Fast (approximate) search?

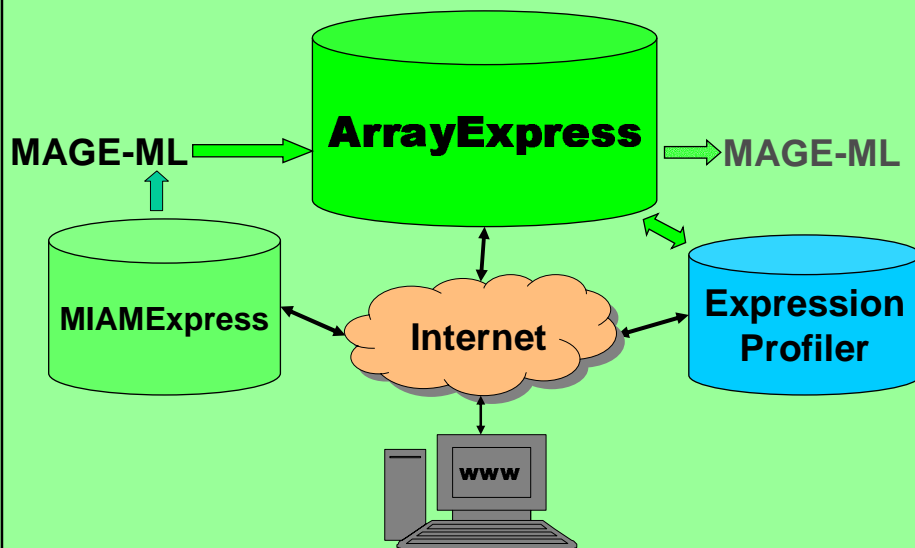
**CCTAGTAG**

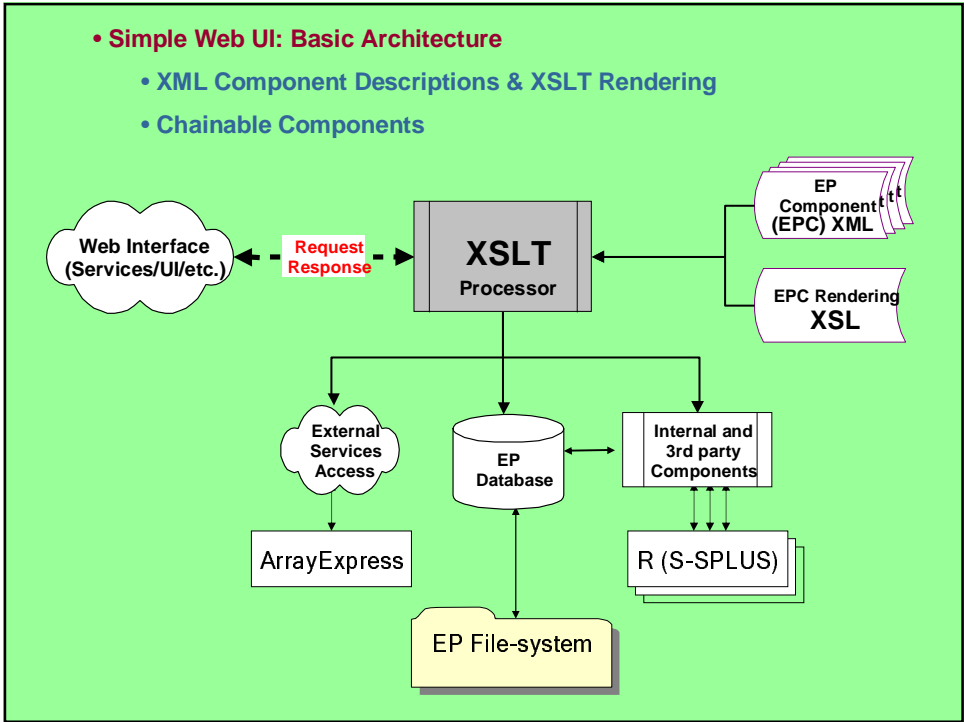
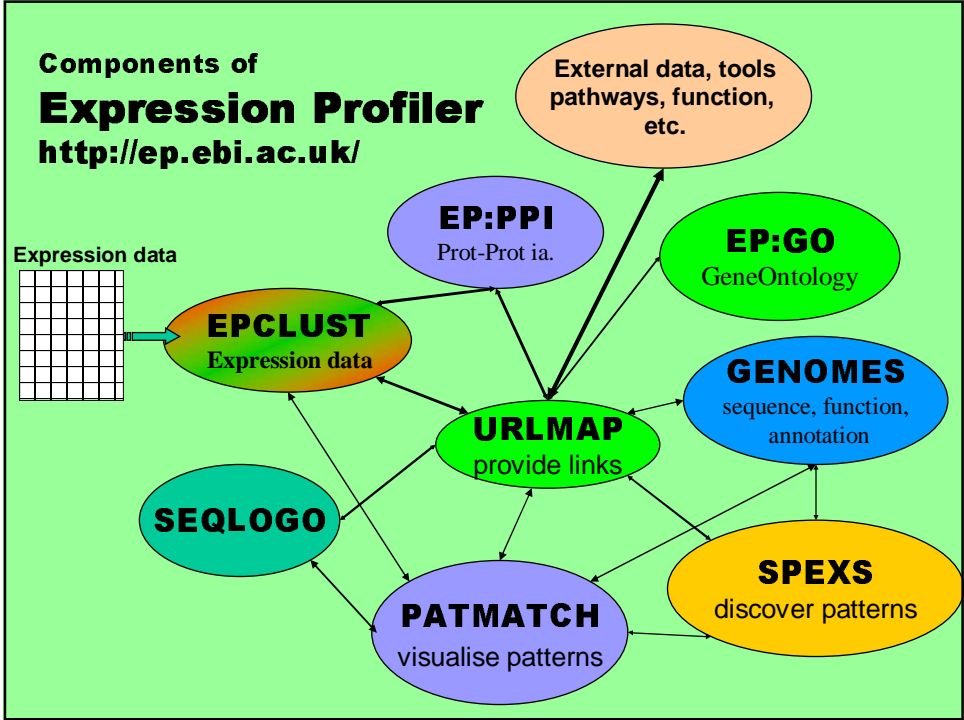


**GTAA..CCT..CCT**

## ArrayExpress (EBI)

a public repository for microarray gene expression data





## SUBSELECT

Expression Profiler  
(component interface)

## Projects started at Tartu

- Database for gene regulation information
- Tools for using that database
- Pattern matching and discovery
- Alternative splicing regulation (3yr EU project)
- (Fast) gene expression data clustering
- Data mining seminar series;
- other DM and ML methods
- ...

## **Gene regulation**

- **Promoter analysis, also in higher eukaryotes**
- **Alternative Splicing data analysis**
- **Genetic networks**
- **Gene expression data analysis**
- **Integration of many different data types**
  - **Protein-protein interactions**
  - **Phenotypes**
- **Metabolic pathways**
- **Signaling pathways**

## **Pattern discovery**

- **Pattern discovery and pattern matching in sequences; sequence algorithms**
- **Regulatory sequence analysis**
- **GPCR receptor bioinformatics**

## Data mining

- **Data mining methods development for bioinformatics**
  - Fast clustering methods
  - Gene networks and regularities
  - Machine learning methods
- **Text Mining**
  - Information extraction
  - categorization
- **Information retrieval, dictionaries**
- **Medical and clinical data handling and storage, population and statistical genetics, pharmacogenetics.**

## Software engineering

- **Database development, software engineering**
  - UML based development and code generation
  - XML based UI-s
- **Expression Profiler**
- **Farm and GRID computing**

## Acknowledgements



**Alvis Brazma**  
**Misha Kapushesky**  
**+ the EBI microarray team**

**Frank Holstege**, and **Patrick Kemmeren**, UMC Utrecht

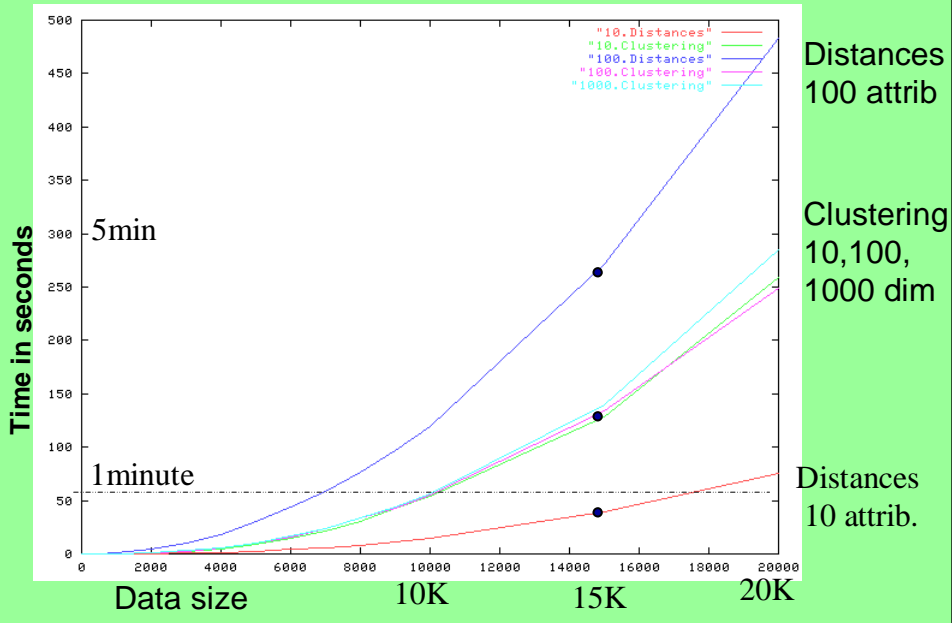
**Mike Croning**, **Steffen Möller**  
(ex EBI)

**Esko Ukkonen**, U. of Helsinki  
**Inge Jonassen**, Bergen U.

**Meelis Kull**  
**Hedi Peterson**  
**Ireen Meho**  
...  
**+ ~ 10 new faces**



## Running time for hierarchical clustering

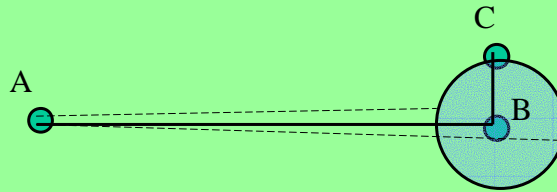


## Limits of standard clustering

- Hierarchical clustering is (very) good for visualization (first impression) and browsing
- Speed for modern data sets remains relatively slow (minutes or even hours)
- ArrayExpress database needs some faster analytical tools
- Hard to predict number of clusters (=>Unsupervised)

## Approximate distances

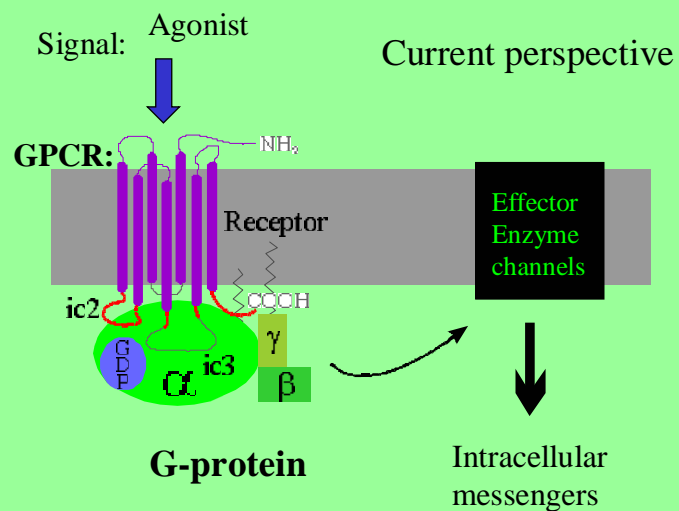
- Triangle inequality for metrics



$d(A,B)$  and  $d(B,C)$  allow us to estimate  $d(A,C)$  within certain limits

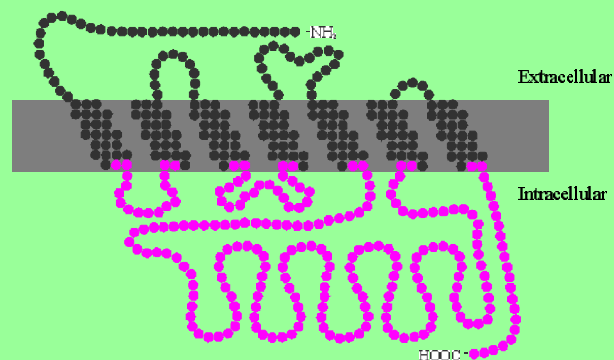
$$|d(A,B) - d(B,C)| \leq d(A,C) \leq d(A,B) + d(B,C)$$

## GPCR coupling

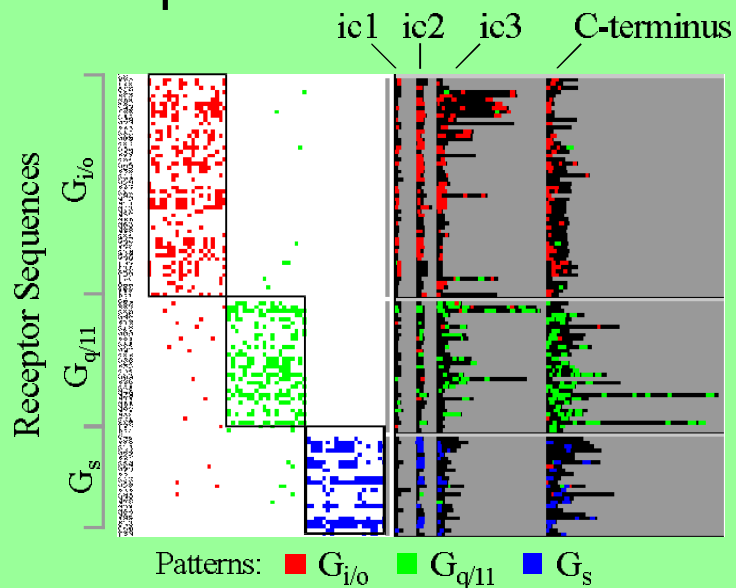


## Our Computational Approach

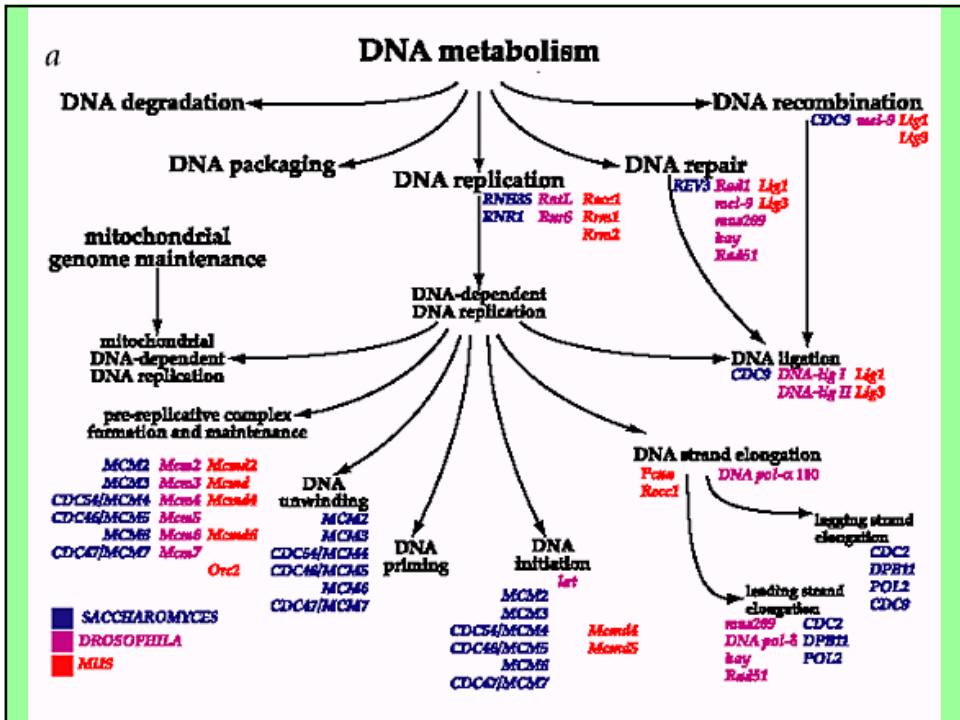
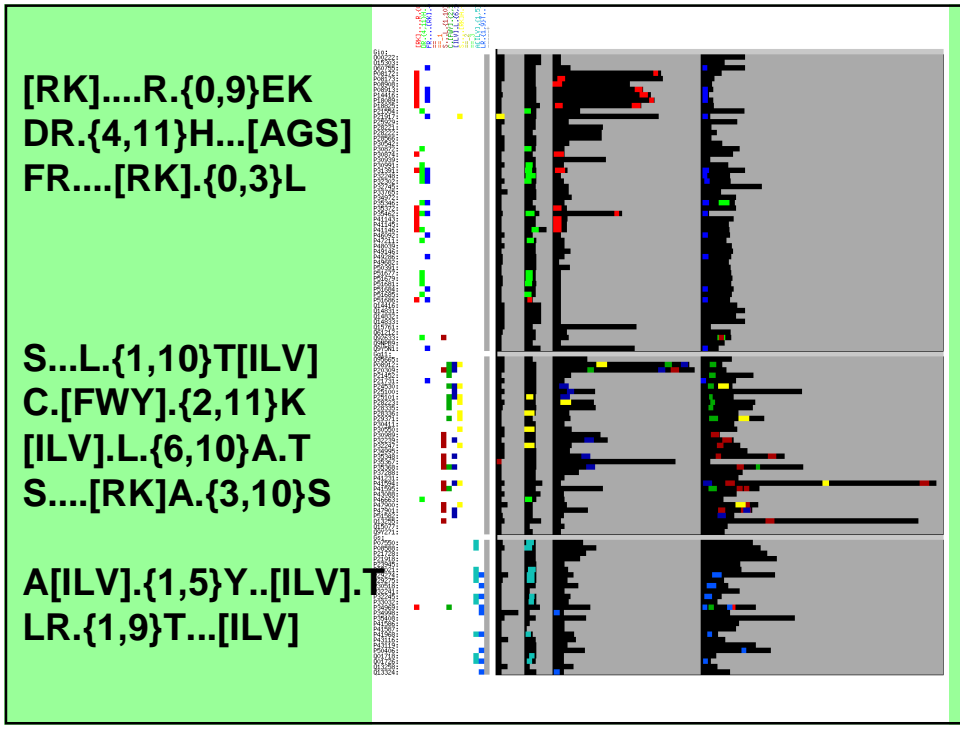
- Using a new membrane topology prediction algorithm (designed specifically for GPCRs), we constrained our pattern search to the intracellular domains of  $\approx 100$  receptor sequences with well-characterised, and non-promiscuous coupling (split into  $G_s$ ,  $G_{i/o}$  and  $G_{q/11}$ )



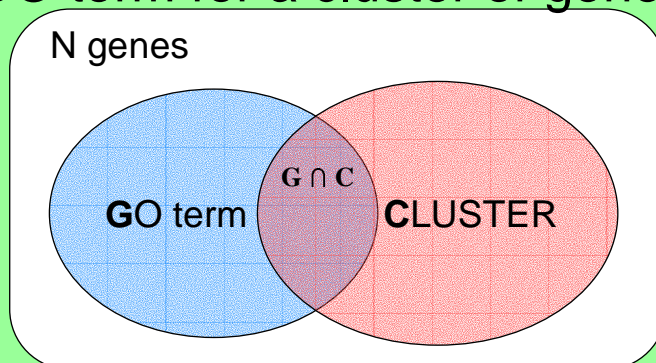
## Receptor Match Positions



Croning, Vilo, Möller, *ISMB 2001*



## Determine the significance of GO term for a cluster of genes



A:  $|G \cap C| / \min(|G|, |C|)$

B:  $P(\text{choose } |C| \text{ from } N \text{ with } |G|, \text{ observe } |G \cap C|+)$

## Annotation of clusters

[GO:0042254](#) <U:L> Process: ribosome biogenesis and assembly (+2:15) (depth=7) [sgd:2:187]

GO:0042254: 47 from cluster (size 98) vs 187 in this class (including subclasses)

[GO:0006364](#) <U:L> Process: rRNA processing (+3:3) (depth=8) [sgd:50:126]

GO:0006364: 35 from cluster (size 98) vs 126 in this class (including subclasses)

[GO:0006360](#) <U:L> Process: transcription from Pol I promoter (+6:14) (depth=8)

[sgd:23:155]

GO:0006360: 38 from cluster (size 98) vs 155 in this class (including subclasses)

[GO:0005730](#) <U:L> Component: nucleolus (+10:17) (depth=6) [sgd:154:210]

GO:0005730: 45 from cluster (size 98) vs 210 in this class (including subclasses)

[GO:0030515](#) <U:L> Function: snoRNA binding (depth=6) [sgd:23:23]

GO:0030515: 17 from cluster (size 98) vs 23 in this class (including subclasses)

[GO:0030490](#) <U:L> Process: processing of 20S pre-rRNA (depth=9) [sgd:33:33]

GO:0030490: 18 from cluster (size 98) vs 33 in this class (including subclasses)

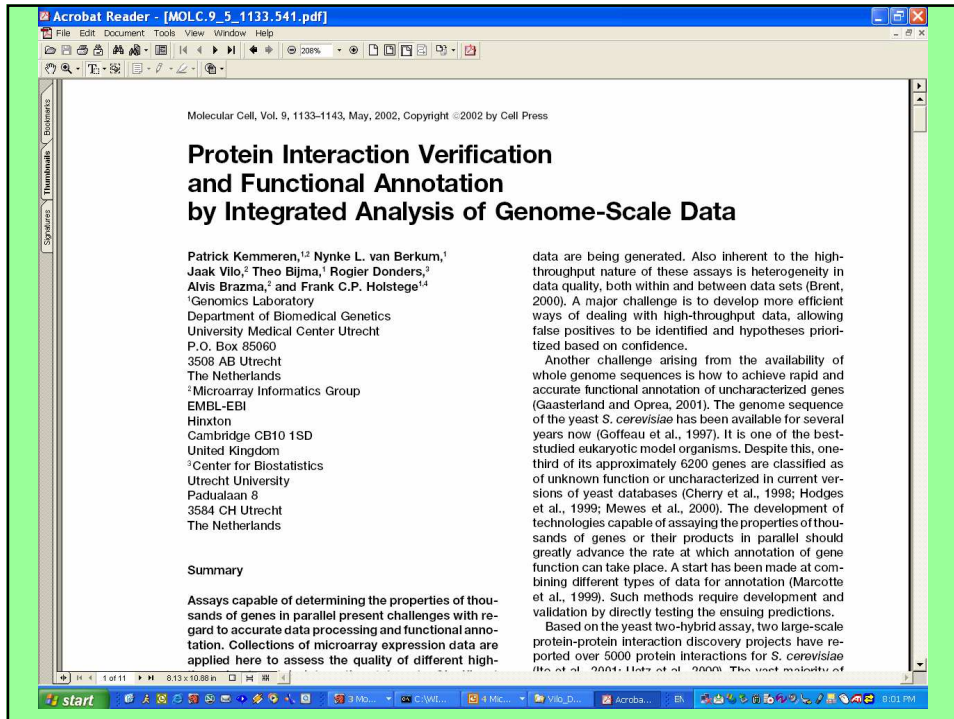
[GO:0005732](#) <U:L> Component: small nucleolar ribonucleoprotein complex (depth=6)

[sgd:30:30]

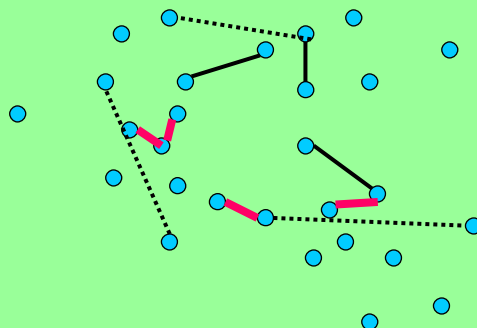
GO:0005732: 16 from cluster (size 98) vs 30 in this class (including subclasses)

[GO:0006396](#) <U:L> Process: RNA processing (+7:52) (depth=7) [sgd:7:370]

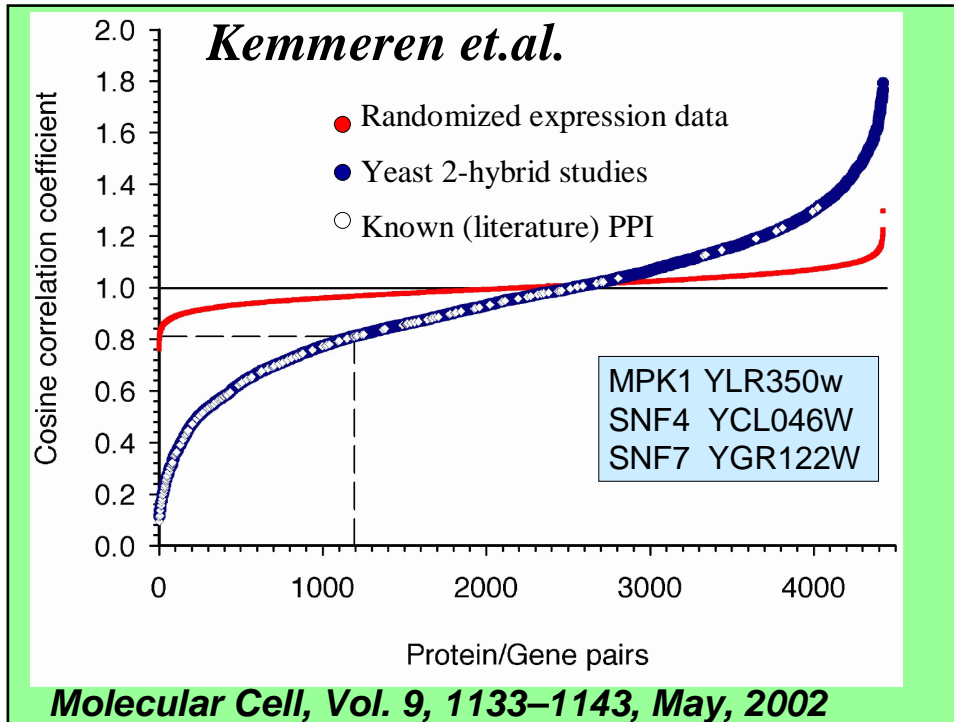
GO:0006396: 40 from cluster (size 98) vs 370 in this class (including subclasses)



## Protein-protein interactions: which to trust more?



**Answer: Use the distance measure alone**



## Results from PPI & expression

- Confidence in 973 out of 5342 putative two-hybrid interactions from *S. cerevisiae* is increased.
- Besides verification, integration of expression and interaction data is employed to provide functional annotation for over 300 previously uncharacterized genes.
- The robustness of these approaches is demonstrated by experiments that test the in silico predictions made.
- This study shows how integration improves the utility of different types of functional genomic data and how well this contributes to functional annotation.



Start summarizing...