

# Gene Regulation: Bioinformatic aspects

Jaak Vilo

CS theory days, Koke, 4.2.04

## Topics

- Biological background
- Computational methods/challenges
- Current projects

### 300+ Cell types

+Brain has ~10,000

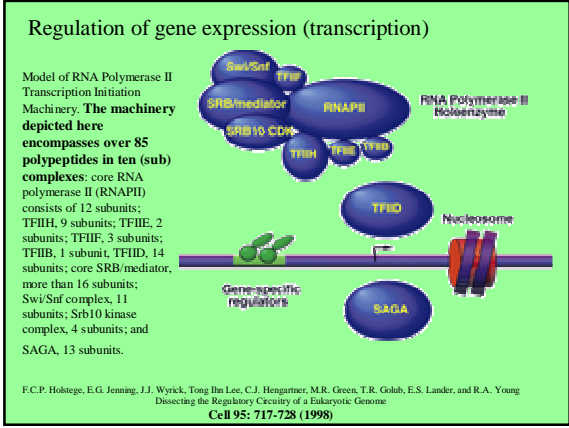
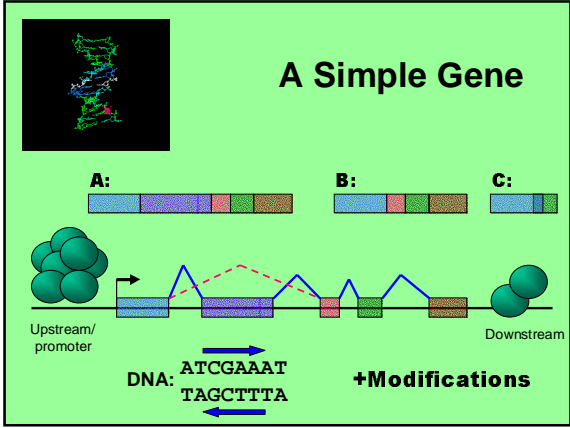
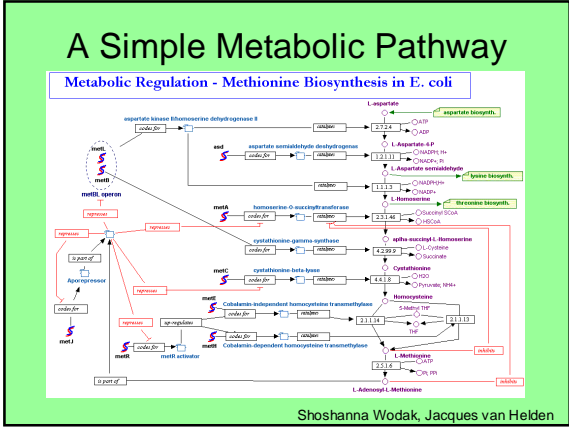
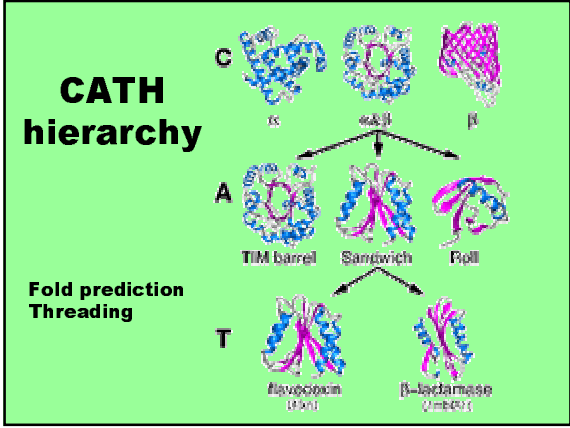
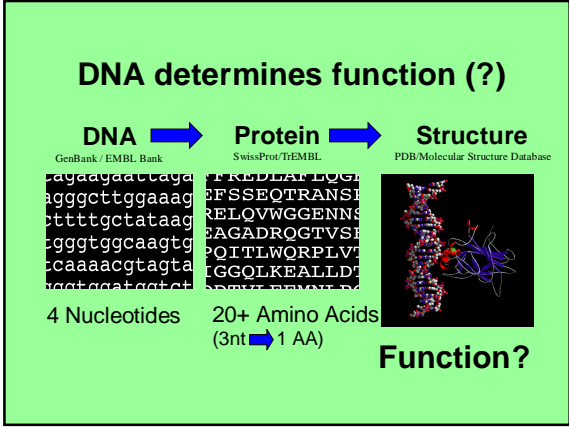
The diagram illustrates the differentiation of a zygote into various cell types. It shows the zygote developing into a blastocyst and then a gastrula. From the gastrula, three germ layers emerge: the ectoderm (outer layer) which gives rise to neurons and pigment cells; the mesoderm (middle layer) which gives rise to muscle cells, blood cells, and cells of the kidney; and the endoderm (inner layer) which gives rise to lung cells, liver cells, and pancreatic cells. A detailed view of a heart cell shows ion channels and G-protein coupled receptors on its cell membrane. A myofibril is also shown, composed of actin and myosin filaments.

David S. Goodsell  
<http://www.scripps.edu/pub/goodsell/>

## Central dogma

TTAAGCTCCG TAGCA DNA  
 ↓  
UUAAGCTCCG TAGCA mRNA  
 ↓  
Leu Ser Ser Val Ala vauk

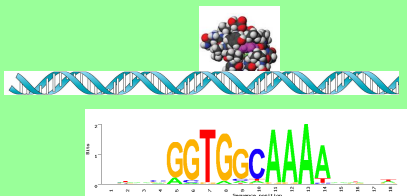
<p>Level 0</p> <p>Level 1</p> <p>Level 2</p> <p>Level 3</p> <p>Level 4</p> <p>Level 5</p> <p>Level 6</p>	<p>ATCGCTGAATTCCAATGTG</p> <p>short region of DNA double helix</p> <p>"beads-on-a-string" form of chromatin</p> <p>30-nm chromatin fiber of packed nucleosomes</p> <p>section of chromosome in an extended form</p> <p>condensed section of metaphase chromosome</p> <p>entire metaphase chromosome</p>	<p>A eukaryotic genome can be thought of as six Levels of DNA structure.</p> <p>The loops at Level 4 range from 0.5kb to 100kb in length.</p> <p>If these loops were stabilized then the genes inside the loop would not be expressed.</p>
--	---	--



### Gene regulation

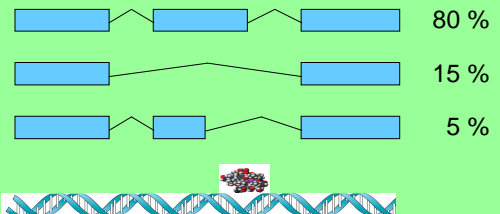
- Determines
  - the development (from embryo)
  - cell types
  - processes of the cell
  - response to the environment
  - ...
- Regulation happens at different levels

## Regulation by binding to DNA/RNA



$$4^6 = 4096, 4^8 = 65.000$$

## Regulation of splicing



Valgu seondumine võib mõjutada splaiisingut

## Regulation of Alternative Splicing

- Which splice variants in which cells?
- Are there cell type specific splicing regulators and signals in DNA/RNA?
- Find genes that have an exon switched on specifically in tissue X
- Is there a common signal for all such exons or splicing events?

## Tissue specific alternative splicing

EST-tehnoloogial baseeruvad andmed (Meelis Kull)

		T1	T2	T3	T4	T5	T6	T7	T8	T9	T10	sum
Geen 1	V1	1	1	0	0	3	0	1	2	1	6	15
	V2	5	0	1	3	3	2	2	9	4	9	38
	V3	3	0	1	2	1	0	0	1	0	1	9
	V4	1	0	0	1	0	0	0	0	1	0	3
	V5	0	0	0	0	1	0	0	0	0	0	3
Geen 2	V1	8	1	3	4	1	1	2	11	3	12	46
	V2	3	0	3	0	0	0	2	7	0	4	19
	V3	0	0	0	0	1	1	1	0	0	1	4
	V4	2	0	0	0	0	0	1	2	1	2	8
	V5	0	0	0	0	0	0	0	1	0	1	2
	V6	0	0	1	0	0	0	1	0	0	0	2
Geen 3	V1	16	1	3	5	2	4	3	17	7	18	76
	V2	7	0	1	2	0	2	2	6	4	8	32
	V3	1	0	0	0	0	1	2	1	1	1	7

## How to study the gene regulation with computational methods?

- What data is available?
- How to combine them meaningfully?
- Algorithms (is the analysis feasible)?
- Actual analysis
- Interpret the results

## Core data (Static)

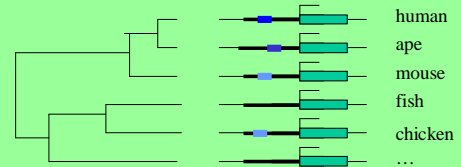
- **DNA sequence(s)**
  - Genes
  - Protein sequences
  - Relation to other species
  - Protein structure (???)
- Partial knowledge about function
  - how to capture this formally?



## Upstream vs genomic random

## Phylogenetic footprinting

Study the same gene in many species

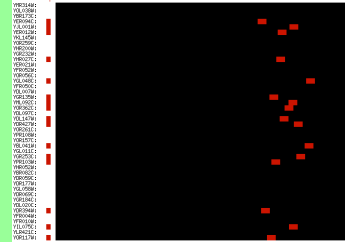


If preserved during evolution then must be important for something!!!

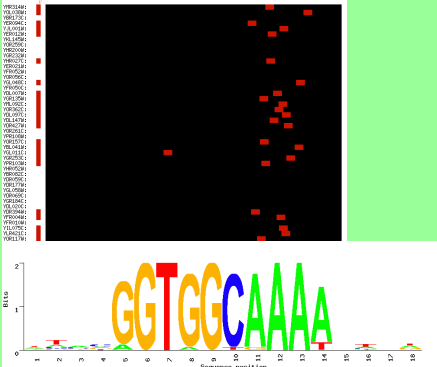
Similar function or role  
→ same regulation?

- This may or may not be true
- How do we actually know that they are behaving similarly?
- Different regulation mechanisms may achieve the same effect

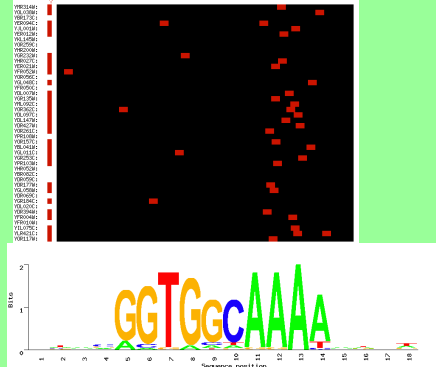
### Proteasome: GGTGGCAAA

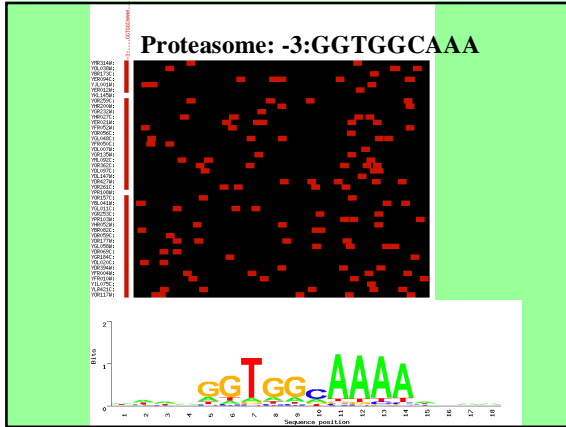


### Proteasome: -1:GGTGGCAAA



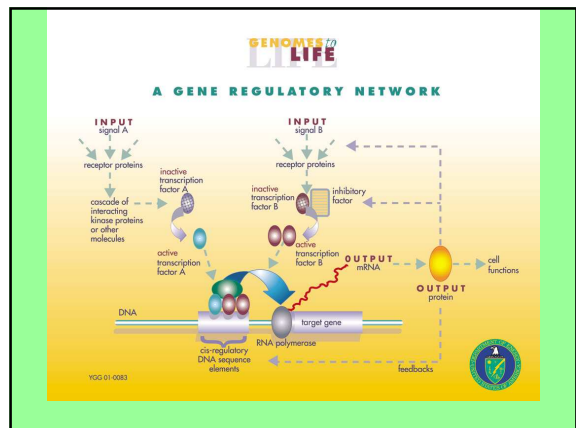
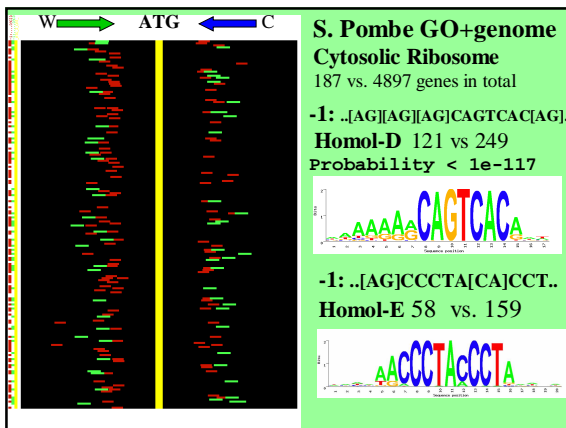
### Proteasome: -2:GGTGGCAAA





## Proteasome movie

- [Movies\proteasome.wmv](#)

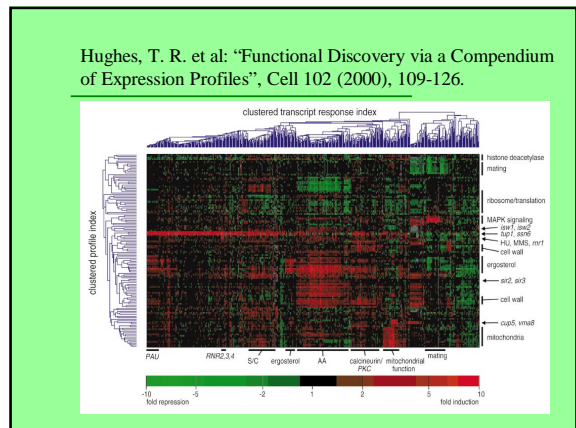
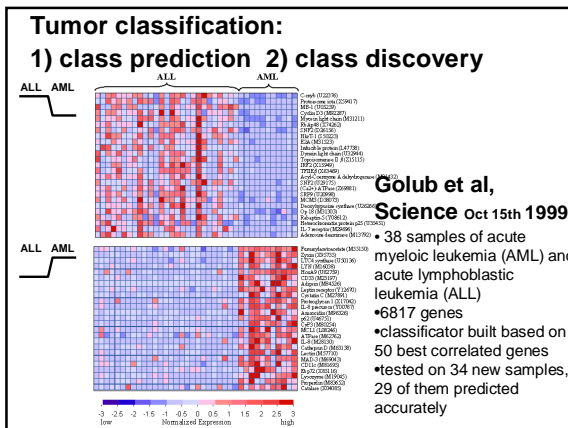
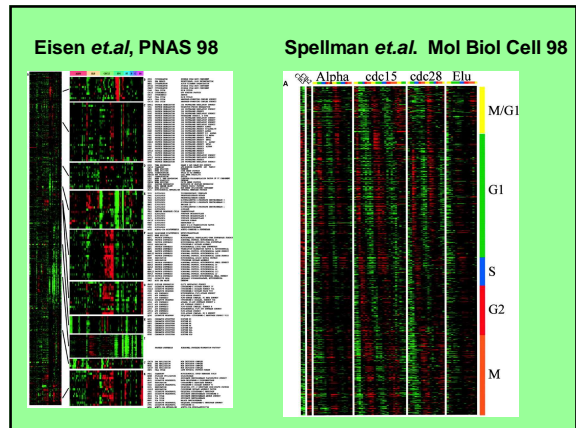
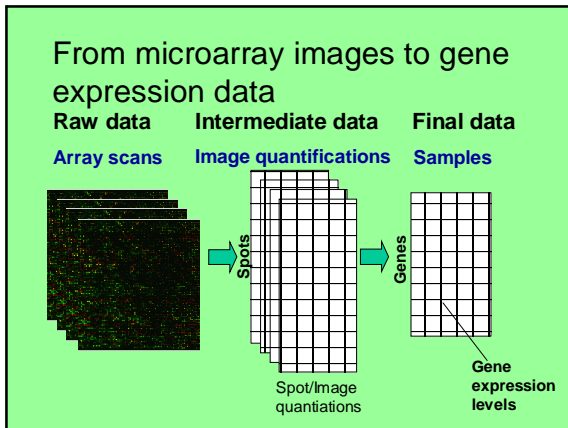
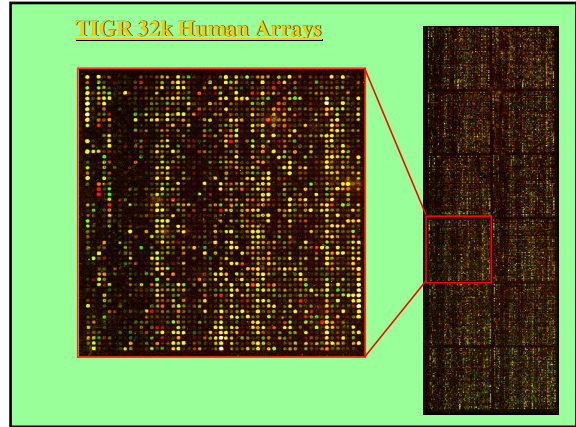
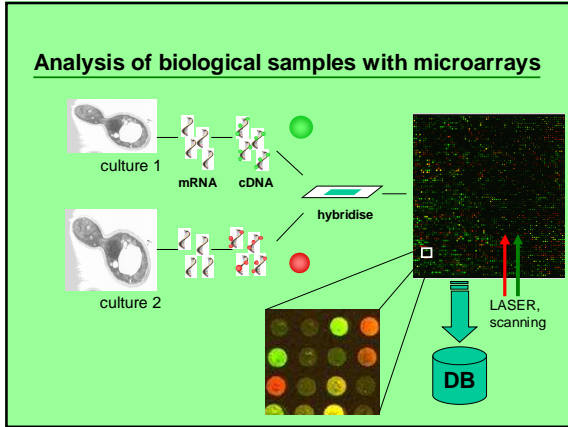


## Dynamics?

- Which genes regulate others
- When and how genes are 'switched on or off?'
- What is the global relationship between genes
- How to model the gene regulation?
  - Continuous stochastic processes responding to the external stimuli

## Experimental data?

- What data can we start with?
- What is known or hypothesised so far?
- Can one test the new hypotheses in practice?

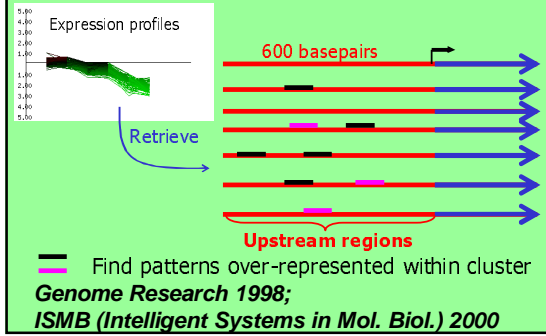




## Gene expression data

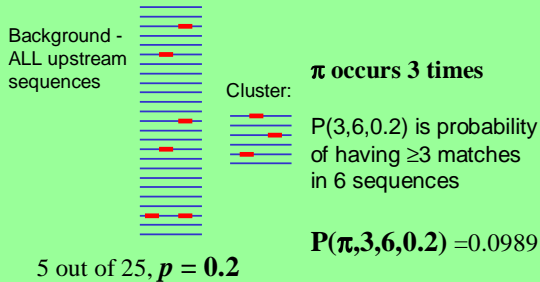
- Snapshots in time to various stimuli, conditions, tissues, time,
- Approximate information about the level of gene expression (RNA transcripts)
- Limited granularity of time
- Limited accuracy
- Data size is large => need fast methods
  - Algorithm: Meelis Kull and J.V.

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



## Pattern selection criteria

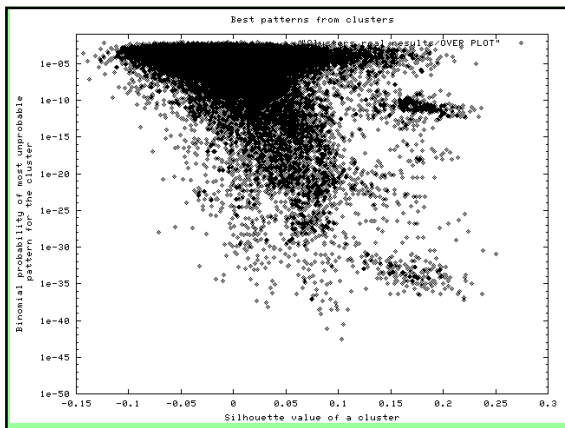
### Binomial distribution



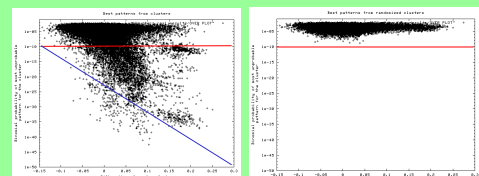
## The most unprobable pattern from best clusters

Pattern	Probability	Cluster size	Occurrences in cluster	Total nr of occurrences in K-mers	K
AAAATTTT	2.59E-43	96	72	930	69
ACGG	6.41E-39	96	75	1088	50
ACGGGT	5.23E-38	94	52	387	40
CCTGGACTAA	5.43E-38	27	18	23	220
GACGG	7.89E-31	86	40	284	38
TTTGGAAACTACAAAAT	2.08E-29	26	14	18	450
TTC TTGTCAAAAAC	2.08E-29	26	14	18	325
ACATACTATTGTAAT	3.81E-26	22	13	18	289
GATGAGATG	5.60E-26	68	24	83	84
TGT TTATATGATGGA	1.90E-27	24	13	18	220
GATGGATTTCGTCAAAA	5.04E-27	18	12	18	300
TATAATAGAGC	1.51E-26	27	13	18	300
GATTTCTGTCAAA	3.40E-26	20	12	18	700
GATGGATTTCCTG	3.40E-26	20	12	18	875
GGTGGCAA	4.18E-26	40	20	96	188
TTC TTGTCAAAAAGCA	5.10E-26	29	13	18	250
GGAAACTTACAAA	5.10E-26	29	13	18	290
GAACTTACAAAATAAA	7.92E-26	21	12	18	550
TTTGTATATG	1.74E-25	22	12	18	600
ATCAACATACTATTG	3.62E-25	23	12	18	375
ATCAACATACTATTGTA	3.62E-25	23	12	18	625
GAACGGCG	4.47E-25	20	11	13	260
GTTAATTTCGAAC	7.23E-25	24	12	18	400
GGTGGCAAA	3.37E-24	33	14	31	475
ATCTTGT TTATATGGA	7.19E-24	19	11	18	875
TTTGTATATGATGGA	7.19E-24	19	11	18	475
GTGGCAA	1.14E-23	28	18	137	725

Vilo *et al.* ISMB 2000



## Significance of the patterns

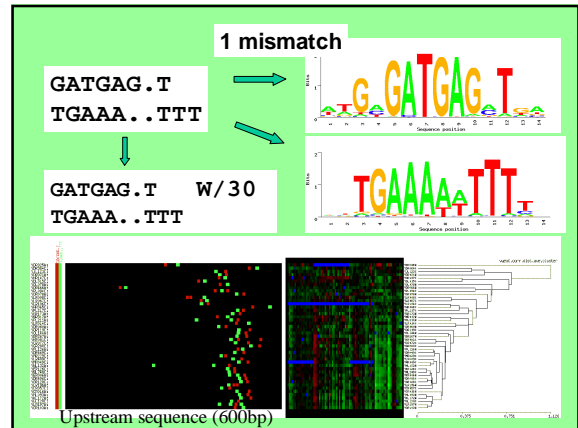
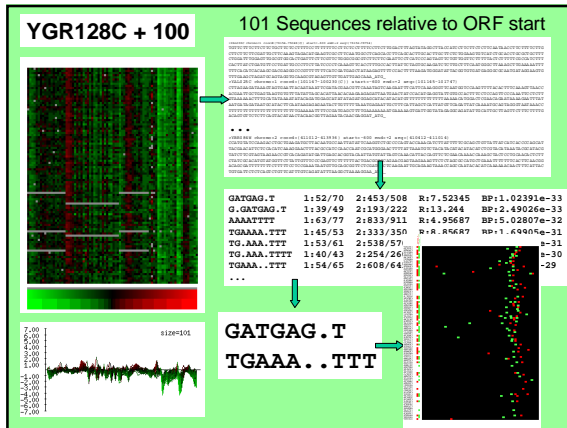


The pattern probability vs. the average silhouette for the cluster

The same for randomised clusters

Vilo *et al.* ISMB 2000





### Problems

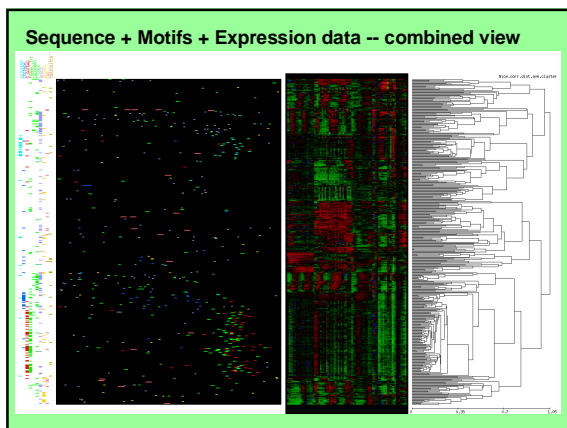
- Many motifs are statistically significant
- Many of them similar to each other

⇒ Summarize meaningfully!  
⇒ Create probabilistic models

- Algorithms: J.V. and Triinu Tasa

### Annotation of clusters

- Map gene sets to GeneOntology categories.
  - GO:0042254 <U> Process: ribosome biogenesis and assembly (+2:15) (depth=7) [sgd:2:187]
  - GO:0042254: 47 from cluster (size 98) vs 157 in this class (including subclasses)
  - GO:0005364 <U> Process: rRNA processing (+3:3) (depth=8) [sgd:50:126]
  - GO:0005364: 35 from cluster (size 98) vs 126 in this class (including subclasses)
  - GO:0005360 <U> Process: transcription from Pol I promoter (+6:14) (depth=8) [sgd:23:155]
  - GO:0005360: 38 from cluster (size 98) vs 155 in this class (including subclasses)
  - GO:005730 <U> Component: nucleolus (+10:17) (depth=6) [sgd:114:210]
  - GO:005730: 45 from cluster (size 98) vs 210 in this class (including subclasses)
  - GO:0030515 <U> Function: snRNA binding (depth=6) [sgd:23:23]
  - GO:0030515: 17 from cluster (size 98) vs 23 in this class (including subclasses)
  - GO:0030490 <U> 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:005732 <U> Component: small nucleolar ribonucleoprotein complex (depth=6) [sgd:30:30]
  - GO:005732: 16 from cluster (size 98) vs 30 in this class (including subclasses)
  - GO:0005395 <U> Process: RNA processing (+7:52) (depth=7) [sgd:1:370]
  - GO:0005395: 40 from cluster (size 98) vs 370 in this class (including subclasses)
- Algorithms: J.V. and Jüri Reimand



### Pattern Discovery

1. Choose the language (formalism) to represent the patterns
2. Choose the rating for patterns, to tell that one pattern is "better" than other
3. Design an algorithm that **finds the best patterns** from the pattern class, **fast**.

### Patterns: AT

```
TGTTCTTTCTTCTTTCATACCTTTTCTTTTTTCC
TTCTCCTTTCCTTCTGACTTTTATAGGCTTACCA
TCCTTCTTCTTCAACCTTCTTACATGCTTCTTC
TTCCATGGCTTCAAAGTAGTTCGTGATCCTTCAAT
GCCTCAGCACCTTCAGCACTTGCCTTCCTCTGAA
GTGCTGCACCTGCGCTGTCTGCTATGTTGGAGTT
GGCGTGGCACTGTTTCTTCGACATGGCGGGCTTCT
TCGATTCCTCAGTCTCTAGTTCTGTTGCTTTT
CTCTGATCATGCTCTTTCACTGCTGCTTCCCTG
TGCCCTATCTATATCTCAAAAGTTCACCTTGCCT
TTCCAAGATCTCTCATATGGGCTTAAAGCCGTAC
TTTTTCACTCGATAGCTTAAAGATTTTCCACTTTA
GCTGCTGGCTGGCTTATATACGGTGTGAGGGCGC
TTGAAAGATTTTTCTCACAGCGACGAGGGCCCG
AGTGTTTGAGCTAGATCAGTAGGTGCAGCGTAGAGT
CTTAGAAGATAAAGTAGTATACATAGATTCGATC
```

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

```
TGTTCTTTCTTCTTTCATACCTTTTCTTTTTTCC
TTCTCCTTTCATTTCTGACTTTTATAGGCTTACCA
TCCTTCTTCTTCAAAAGTAGTTCGTGATCCTTCAAT
TTCCATGGCTTCAAAGTAGTTCGTGATCCTTCAAT
GCCTCAGCACCTTCAGCACTTGCCTTCCTCTGAA
GTGCTGCACCTGCGCTGTCTGCTATGTTGGAGTT
GGCGTGGCACTGTTTCTTCGACATGGCGGGCTTCT
TCGATTCCTCAGTCTCTAGTTCTGTTGCTTTT
CTCTGATCATGCTCTTTCACTGCTGCTTCCCTG
TGCCCTATCTATATCTCAAAAGTTCACCTTGCCT
TTCCAAGATCTCTCATATGGGCTTAAAGCCGTAC
TTTTTCACTCGATAGCTTAAAGATTTTCCACTTTA
GATCGTGGCTGGCTTATATACGGTGTGATGAGGGCGC
TTGAAAGATTTTTCTCACAGCGACGAGGGCCCG
AGTGTTTGAGCTAGATCAGTAGGTGCAGCGTAGAGT
CTTAGAAGATAAAGTAGTATACATAGATTCGATC
```

### SPEXS - Sequence Pattern EXhaustive Search Jaak Vilo, 1998

- **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

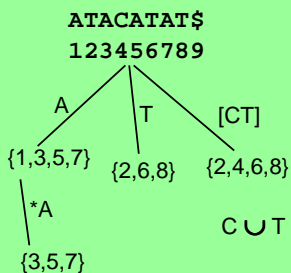
### Regular patterns

- 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.....GCA**

### 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

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

## How to improve?

- **Simple vs complex patterns/profiles**
  - What is the best representation?
  - What is the best algorithmic approach
- Can we prove/disprove expression data clustering methods or distance measures by systematic promoter analysis?
- Lots of computations to perform...
- Tools for non-algorithm persons - how to maintain the simplicity vs desired results vs computational complexity

## Implant $k,d$ -patterns

(The Challenge Problem, P.Pevzner, 2000)

Length  $k=15$

TGATTTCTTCGACAT

$d=4$ , nr of changed characters

TGTTATCTTGGAGAT

TGAATTGTTCCACAC

Such motifs can differ in up to 8 positions out of 15!

```

TGATTTCTTCGACATCGTTTTCCTTTTFFCC
TTCTCCTTTGATTTCCGACTTTTAAATAGGCTTACA
TCCATCTCTCTGATAGACCTTCTTACATGCTTCTC
TTGATTTGCTTAAAGTAGTTGGGATCATCTTCAAT
GGCTCAGCACCTTCAGGACTTGCCTTCACTTCTGGAA
GTCTGACCTGGGCTTTCTTCTGAAAGATTTGGAGT
GGCTGGCAGTATTCTTCGACATGGGCGGCTTCTT
TGGATTTCCATGAGCTTATAGTTCGTGTGTTTFTT
CTGAAATGATGGTCACTTTCAGTACTTGAATGCTGG
TCCCTATATATCATCTAAAGTTCACTTGGCCTAC
TTCCAGATCTCTCATCATATAAGGCTTAAAGCCATC
TTTTTCACTGAGAGCTATAAGAGTTTTCACCTTTTA
GATCTGCTGGCTTATATACCTTGAATAGGCGCC
    
```

## Approximate all against all

- Assume at least one perfect occurrence exists
- Only  $O(kn)$  different substrings of length  $k$
- Match all of them approximately
- Find the one that has most significant nr of approximate occurrences
- Trie-index the sequences first, then search
- Algorithm: J.V. and **Hendrik Nigul**

## Gene regulation is affected by

- DNA/RNA sequence
  - signals along that sequence
  - DNA structure and state
- State of the cell
  - i.e. all the other molecules and
- Environment

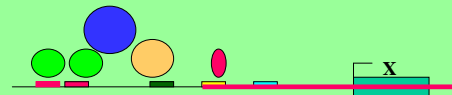
## Binding sites: individually and in combination



Episode rules: A followed by (C D or D C)  
Asko Tiidumaa

Conservation of distances between sites  
Jelena Zaitseva

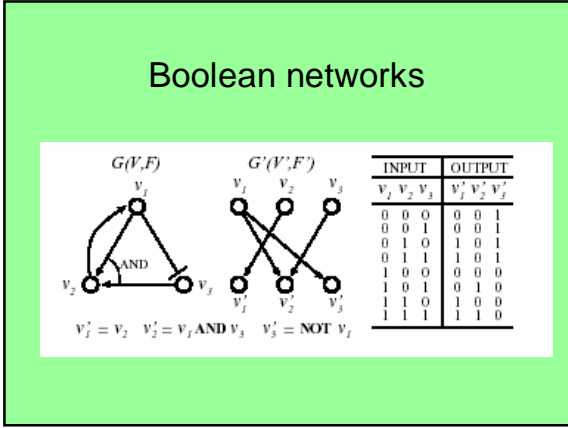
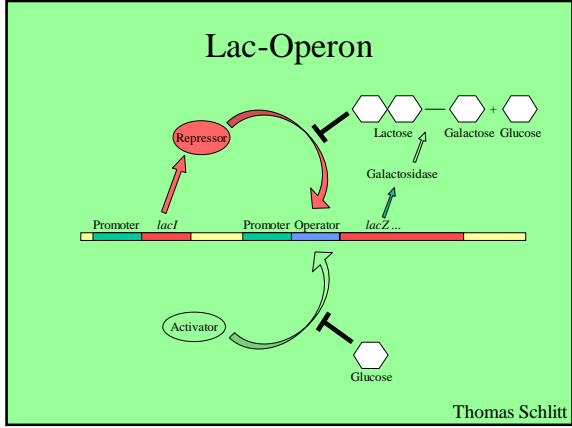
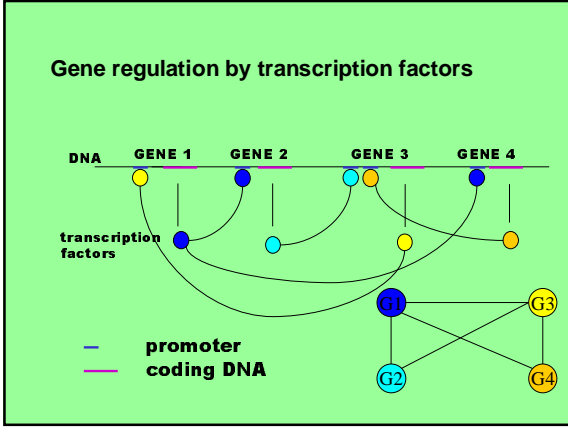
## Goals:



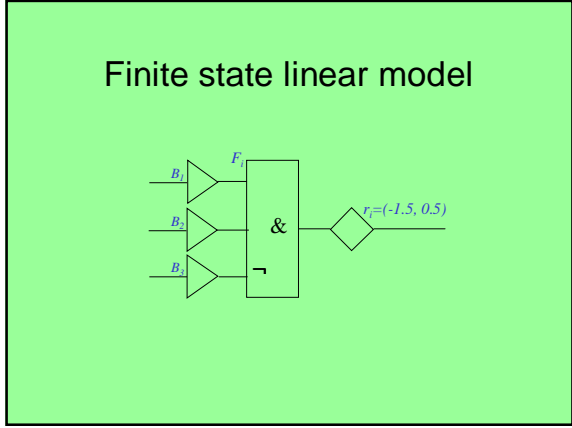
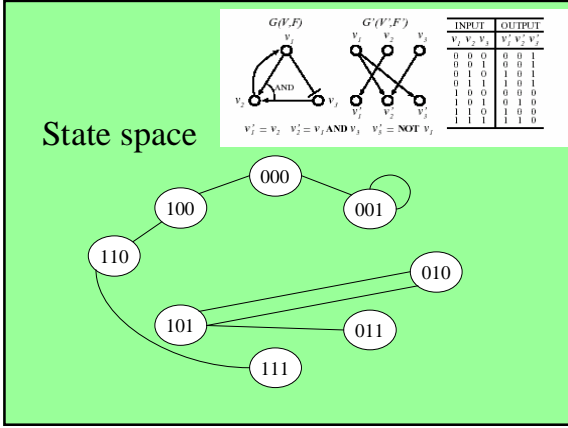
Given the sequence (signals)  
and gene expression levels of other genes:  
predict expression level of gene X

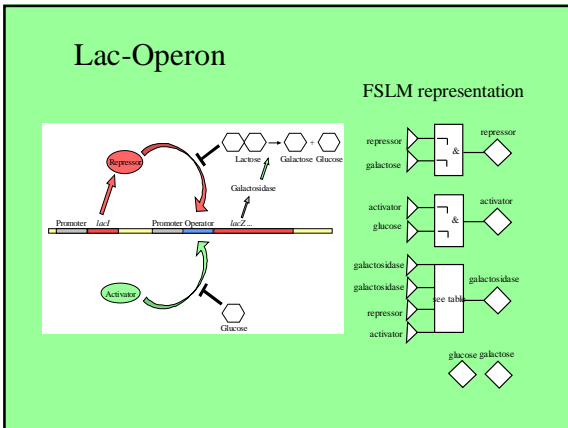
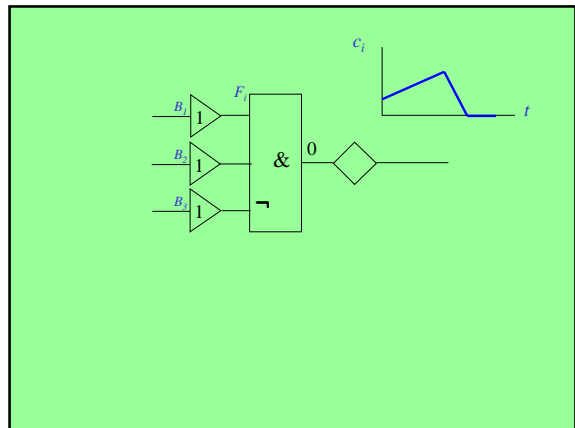
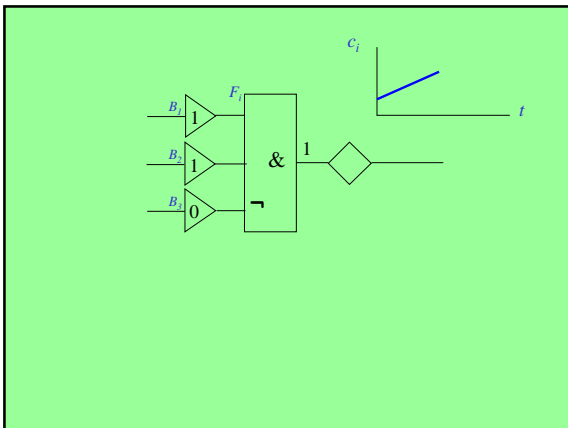
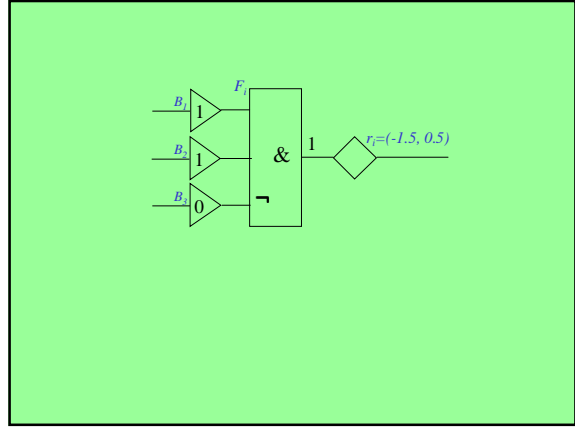
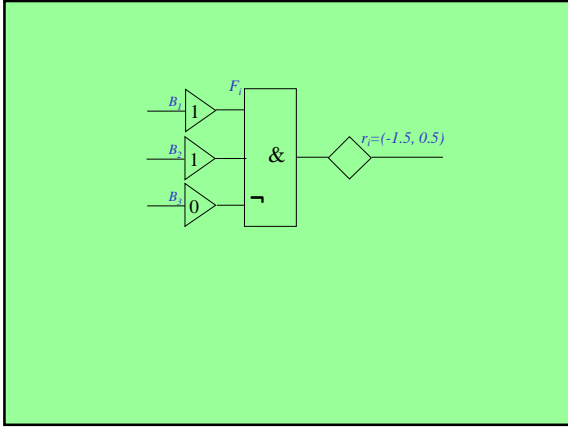
Given the chromosomal sequence  
predict locations of promoters and genes

Predict dependencies between all genes,  
and study gene regulation networks



- ### Synchronous Boolean networks - assumptions in gene network modelling
- Each gene the system (cell) can be in one of **two states** –
    - 'expressed' – 1,
    - 'not expressed' – 0
  - The genes can switch from state to state all simultaneously in **synchronous** manner
  - The next state of each gene is **determined** by previous states of all genes by Boolean functions describing the network





- ### Main related sub-projects
- Clustering – Meelis Kull
  - Motif discovery – Hendrik Nigul, Triinu Tasa, ...
  - Site combinations: Jelena Zaitseva, Asko Tiidumaa
  - Database of Gene Regulation - Hedi Peterson, Eero Raudsepp, ...
  - Annotate sets of genes based on quilt by association – Jüri Reimand
  - Alternative Splicing – Meelis Kull
  - Software development, visualization, GRID, Web Services, etc.