

Teadmiste kaevandamine bioloogilistest andmetest

Jaak Vilo

Teoreetilise arvutiteaduse talvekool,
Arula, 3. veebruar 2003

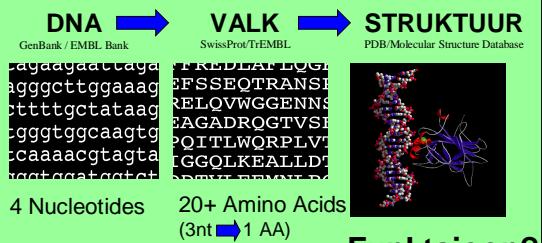
Eesmärgid

- Tutvustada bioloogiliste andmete analüüsiga (kaevandamise) probleemide
- Kas bioinformaatikal on midagi pakkuda teoreetilisele arvutiteadusele?
- Näiteid bio-andmete kaevandamisest

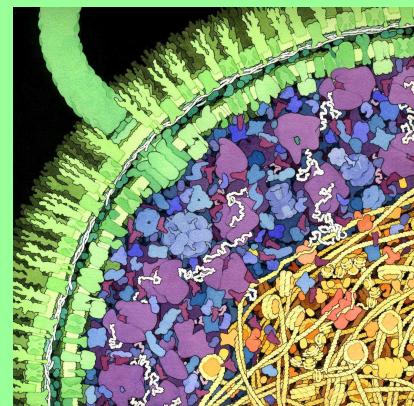
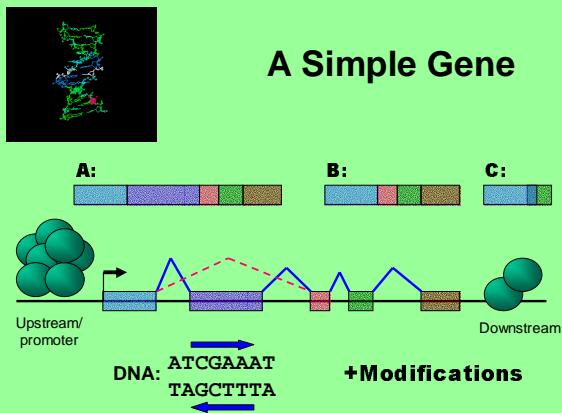
What is Bioinformatics?

- Bioinformatics is the use of information technology to store and analyze genetic information
- Bioinformatic researchers develop and apply computing tools to extract the secrets of the life and death of organisms from the genetic blueprints and molecular structure stored in digital collections

DNA määrab funktsiooni (?)



A Simple Gene



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

Uurimisküsimusi

- Milline on iga geeni ja geeniprodukti (valk, RNA) funktsioon?
- Kuidas täpselt toimuvad bioloogilised protsessid?
 - Transkriptsioon ja translatsioon ja nende reguleerimine
 - Millised valgud interakteeruvad ja miks
 - Metaboolsed rajad ja võrgud
 - Signaali ülekanne rakus ja organismis

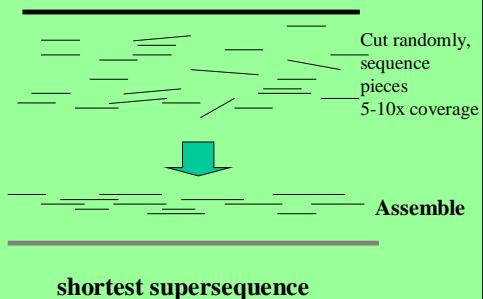
Kus vaja arvutiteadust?

- Andmete kogumise protsessi tugi
- Toorandmetest üldistatud andmete saamine
- Andmete haldamine
- Andmete analüüs
- Teadmiste esitamine
- Modelleerimine ja simulatsioon

Excerpts from curricula descriptions

- New experimental techniques generating mass data are being developed.
- Every new biological research idea requires a specifically tailored, algorithmic approach, raising an abundance of challenging questions in algorithm design, analysis and implementation.
- The demands of biosequence analysis require advances in many classical fields of algorithmics. Our recent work concentrates on suffix trees and related index structures, on efficient pattern matching in strings and trees, and on a new, algebraic style of dynamic programming.

DNA (shotgun) sequencing



Data Mining Andmete kaevandamine

- Data Mining ja Knowledge Discovery from Databases (KDD) on uued arvutiteaduse alad, palju ühist statistikaga
- Eesmärgiks on väga suuret andmekogude analüüs
- Otsi (lihtsaid) reegleid ja seoseid mis kehtivad (vähemalt teatud osas) andmetes
- Statistika, masinõppimine, andmebaasiteooria, + domain knowledge

Bio-andmete kaevandamine

1. Millised bio-andmed on olemas ja kuidas neid võiks kasutada
2. Andmete kogumine, puhastamine, ettevalmistamine ja ühendamine
3. Analüüs arvutil (peamine algoritm)
4. Tulemuste esitamine, reeglite tuvastamine, visualiseerimine
5. Tulemuste tõlgendamine

```

TGTTCCTTCTTCTTCATACATCCTTTCTTTTCC
TCCTCCTTCTCATTTCTGACTTTAAATAAGGCTTACCA
TCCTCTCTCTTCATAAACCTCTTACATTGCTTCTTC
TTCGATTGCTCAAAGTAGTTGTAATCATCTTCAT
GCCTCAGCACCTTCAGCACTTGCACCTTCATTCTGGAA
GTGCTCACCTGCGCTGCTTGTAATGGATTGGAGTT
GGCGTGGCACTGATTTCTTCGACATGGCGGCGTCTTCT
TCGAATTCCATCAGTCTCATAGTTCTGGTTCTTT
CTCTGATGATCGTCATCTTCACTGATCTGATGTTCTG
TGCCCTATCTATATCATCTCAAAGTTACCTTGGCACT
TTCCAAGATCTCATTCATAATGGGCTTAAAGCCGTAC
TTTTTCACTCGATGAGCTATAAGAGTTTCCACTTTA
GATCGTGGCTGGGCTTATTTACGGTGTGATGAGGGCGC
TTGAAAAGATTTTCATCTCACAAAGCGACGAGGGCCG
AGTGTGAAAGCTAGATGCACTGAGGTGCAAGCGTAGAGT
CTTAGAAGATAAAAGTAGTGAATTACAATAGATTGATAC

```

A Challenge Problem (P. Pevzner, 2000)

- Insert into every sequence a 15-mer where 4 positions have been randomly changed
- Challenge: discover what was the original sequence inserted
- Why? 4^{15} sequences in total, variants can differ as much as in 8 positions out of 15

Tekstist mustriga otsimine

- Kas **ATGCAGA** esineb tekstis?
- Kas **ATGCAGA** esineb tekstis ligilähedelt (ATCCAGA, ATGCGA)?
- Kas **A[TA].C[CG].{3,7} A[TA].C[CG]** esineb tekstis?
- Kui kaua võtab aega ülaltoodud päringlete vastamine?
- Kuidas eeltöödelda ja indekseerida?
- Suffiksipuud ja massiivid

Patterns: AT

```

TGTTCCTTCTTCTTCATACATCCTTTCTTTTCC
TCCTCCTTCTCATTTCTGACTTTAAATAAGGCTTACCA
TCCTCTCTCTTCATAAACCTCTTACATTGCTTCTTC
TTCGATTGCTCAAAGTAGTTGTAATCATCTTCAT
GCCTCAGCACCTTCAGCACTTGCACCTTCATTCTGGAA
GTGCTCACCTGCGCTGCTTGTAATGGATTGGAGTT
GGCGTGGCACTGATTTCTTCGACATGGCGGCGTCTTCT
TCGAATTCCATCAGTCTCATAGTTCTGGTTCTTT
CTCTGATGATCGTCATCTTCACTGATCTGATGTTCTG
TGCCCTATCTATATCATCTCAAAGTTACCTTGGCACT
TTCCAAGATCTCTCATCTCAAAGCTTAAAGCCGTAC
TTTTTCACTCGATGAGCTATAAGAGTTTCCACTTTA
GATCGTGGCTGGGCTTATTTACGGTGTGATGAGGGCGC
TTGAAAAGATTTTCATCTCACAAAGCGACGAGGGCCG
AGTGTGAAAGCTAGATGCACTGAGGTGCAAGCGTAGAGT
CTTAGAAGATAAAAGTAGTGAATTACAATAGATTGATAC

```

Algoritmid

- Tehniline: loenda kõik etteantud stringides esinevad mustrid ja nende sageused
- Bioloogiline: kogu kokku sarnaselt ekspressoerunud geenide promootorid ja otsi võimalikke transkriptsiooni-faktorite seondumiseks soodsaid mustreid

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

```
TGTTCTTCTTCATTC[ATAC]CCTTTCTTCTTCTTCTC  
TTCTCTTC[ATTCCTGACTTTAAATAAGGCTTACCA  
TCCTCTCTCTCA[ATACCTTCTACATTCCTCTC  
TTCCATGCTTCARAGTAGTTCTGGA[ATACCTCTCAAT  
GCCTCAGCACCTTCAGGACTTGCACTTCATTCCTGGAA  
GTGCTGACCTGCGCTGCTTGTAAATGAAATTGGAGTT  
GGCGTGGCACTGAA[ATCTTCAC[ATGGCGGGCTCTCT  
TCGAATTCCATCACTCC[ATAGTTCTGTTGGTTCTTT  
CTCTGAA[ATCTGCACTTCACGAA[ATGTTCTGCTG  
TGCCCTAATATAC[ATCAAAACTTCACCTTGCACACT  
TTCCAAG[ATCTCTCAATCAT[ATGGCTTAAAGCCGTAC  
TTTTTCACTCGATGAGCTATAAGAGTTTCACTTTTA  
GATCGTGGCTGGCTTATATACCGGTGTC[ATGGGGCG  
TTGAAGAG[ATTTTCATCTACAAAGCCAGGAGGCCCG  
AGTGTGAAAGCTAG[ATCAGTAGGTGCAAGCGTAGAGT  
CTTAGAAG[ATAACTAGTGAATTACATAG[ATTCGATAC
```

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

```
TGTTCTTCTTCATTC[ATAC]CCTTTCTTCTTCTTCTC  
TTCTCTTC[ATATTCCTGACTTTAAATAAGGCTTACCA  
TCCTCTCTCTCA[ATAACTTCTTCAATTCCTCTC  
TTCGATTGCTTCAAAGTAGTCTGCAATCATTCCTCTCAAT  
GCCTCAGCACCTTCAGGACTTGCACTTCATTCCTGGAA  
GTGCTGACCTGCGCTGCTTGTAAATGGATTGGAGTT  
GGCGTGGCACTGATTTCCTC[ATGGCGGGCTCTCT  
TCGAATTCCATCACTCC[ATAGTTCTGTTGGTTCTTT  
CTCTGATGATTC[ATATTCACGATCTGATGTTCCCTG  
TGCCCTATCTATATCATCTCAAAGTTCACTTGGCCACT  
TTCCAAGATCT[ATATCATATGGCTTAAAGCCGTAC  
TTTTTCACTCGATGAGCTATAAGAGTTTCACTTTA  
GATCGTGGCTGGCTTATATACCGGTGTC[ATGGGGCG  
TTGAAGAG[ATTTTCATCTACAAAGCCAGGAGGCCCG  
AGTGTGAAAGCTAG[ATCAGTAGGTGCAAGCGTAGAGT  
CTTAGAAG[ATAACTAGTGAATTACATAG[ATTCGATAC
```

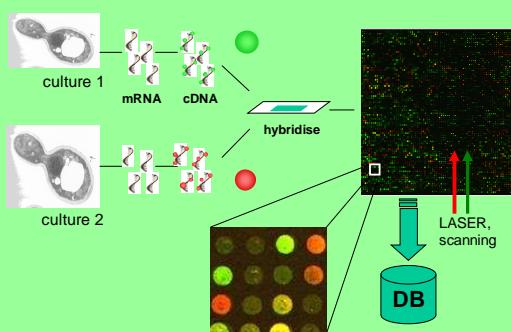
Bioinformaatika algoritmid

- Mitte ainult arvutiressursside küsimus
- Viisid kuidas kombineerida andmeid
- Kas saadav tulemus on bioloogiliselt relevantne
- Kas see edendab meie arusaamist bioloogilistest fenomenidest?

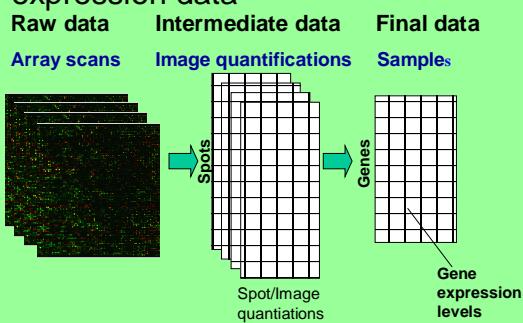
Näidis-meetodid

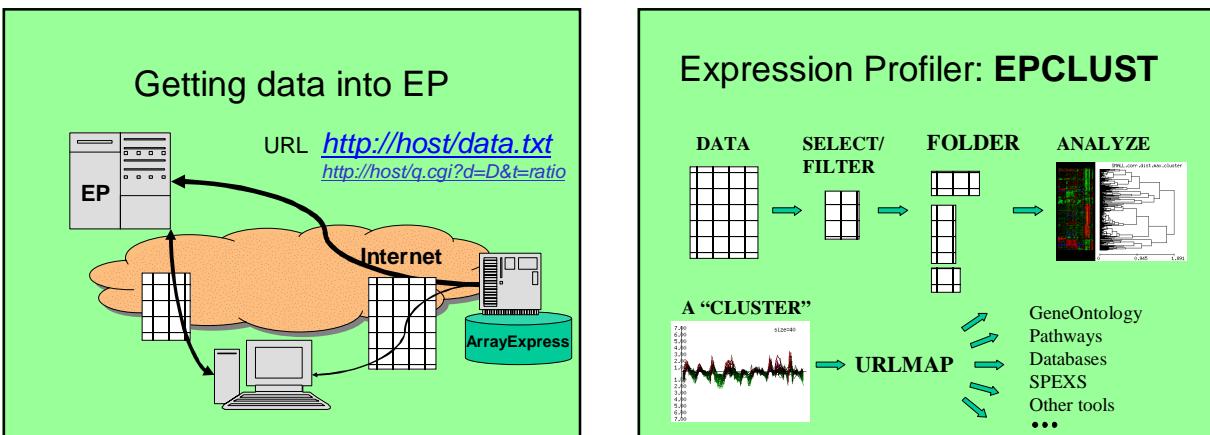
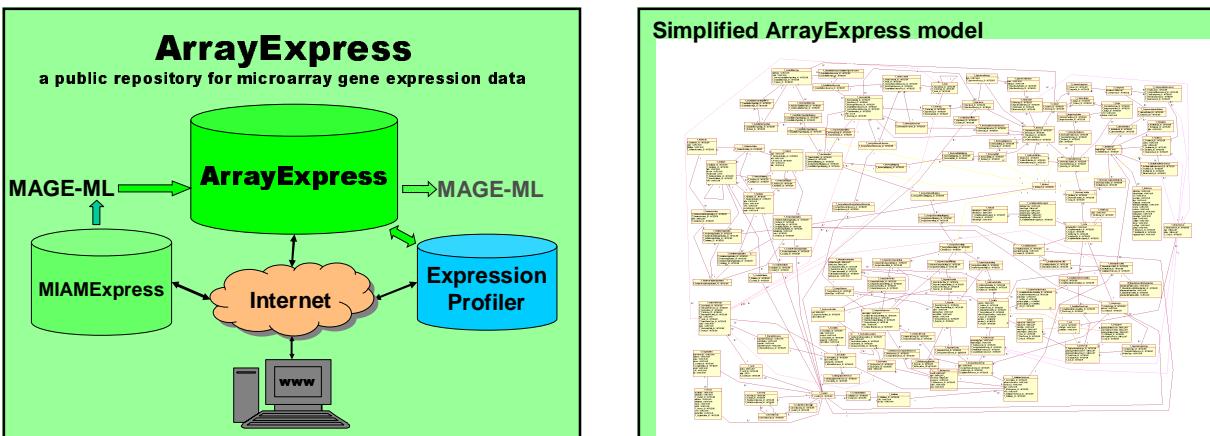
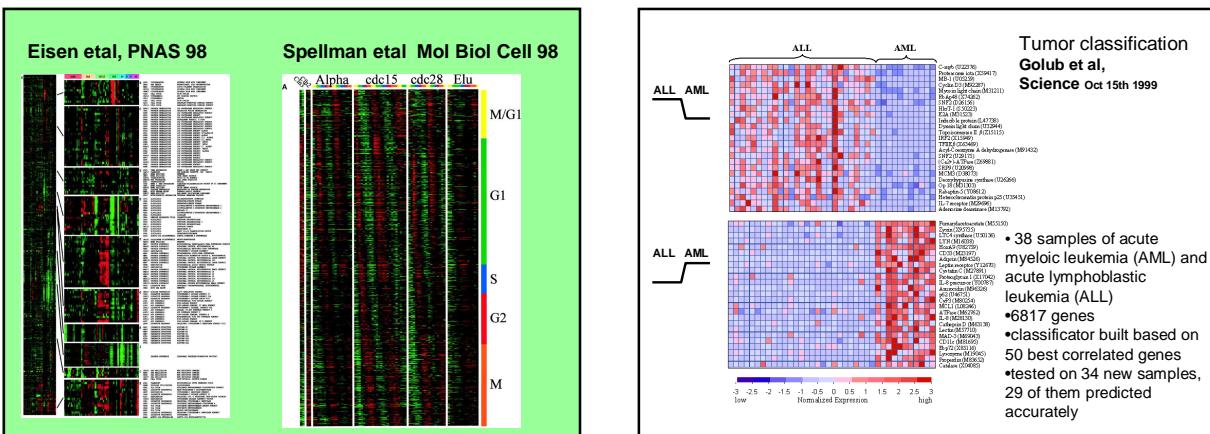
- Geeniekspresiooni analüüs
- Transkriptsioonifaktorite seondumiskohade ennustamine
- Valk-valk interaktsionide analüüs
- G-valk retseptorite ja G-valkude seoste analüüs
- Teadustekstide analüüs

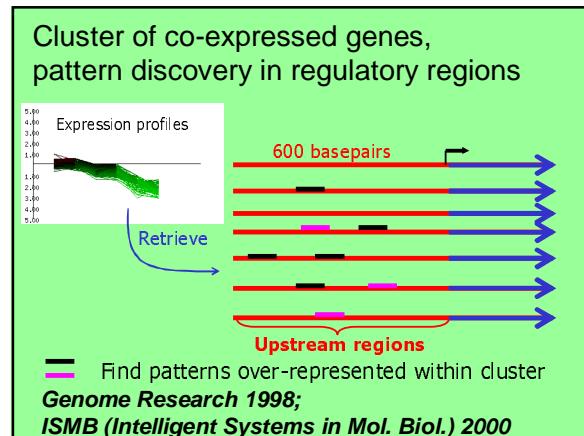
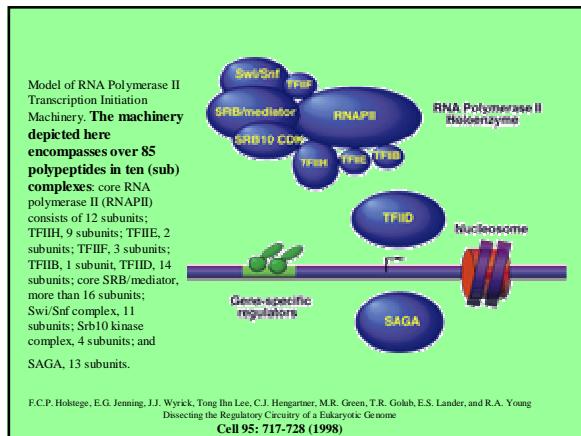
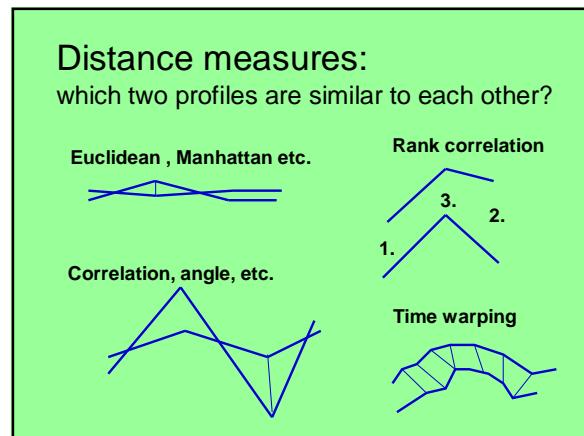
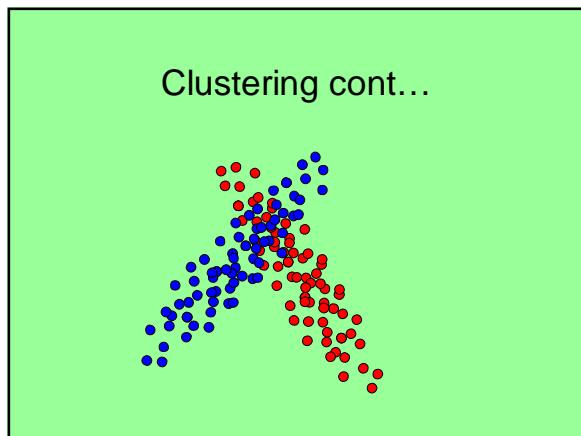
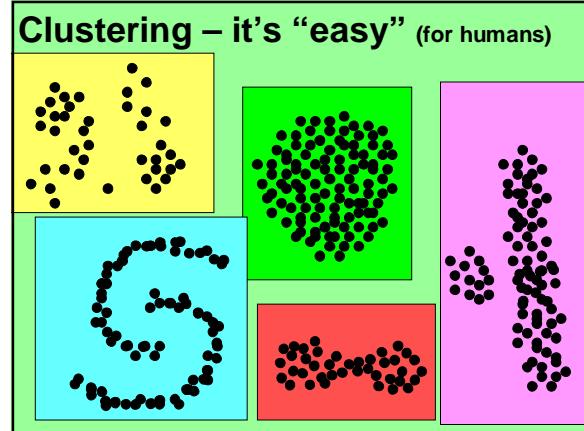
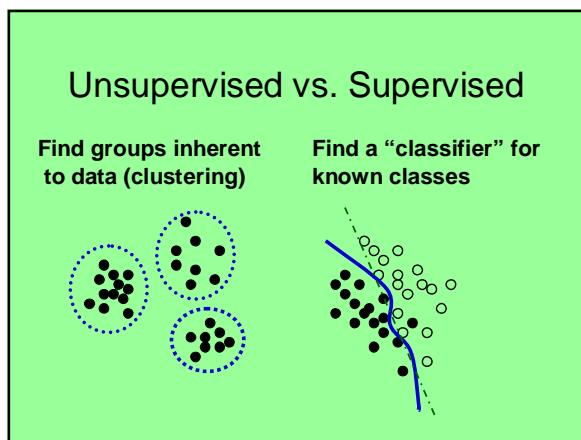
Analysis of biological samples with microarrays

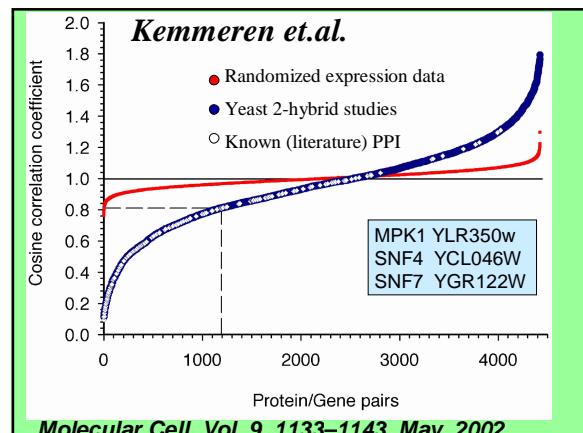
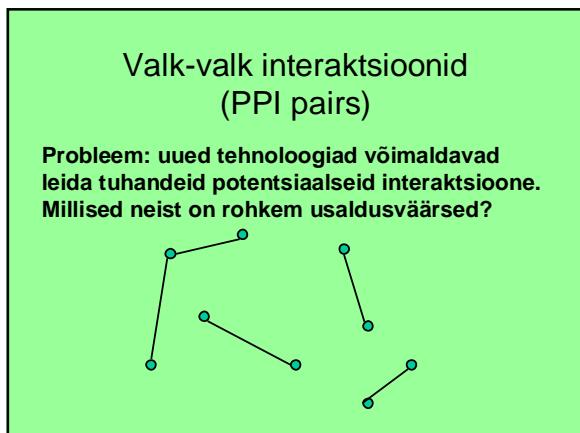
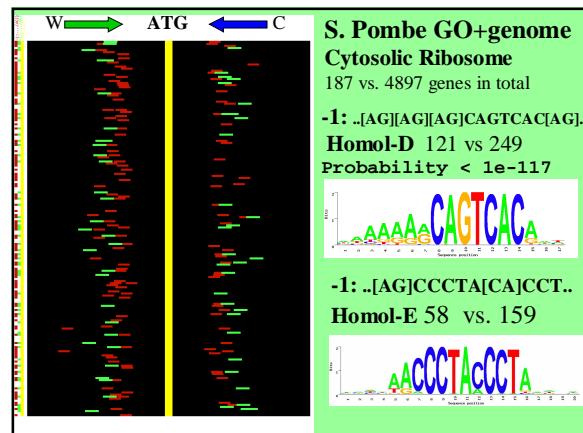
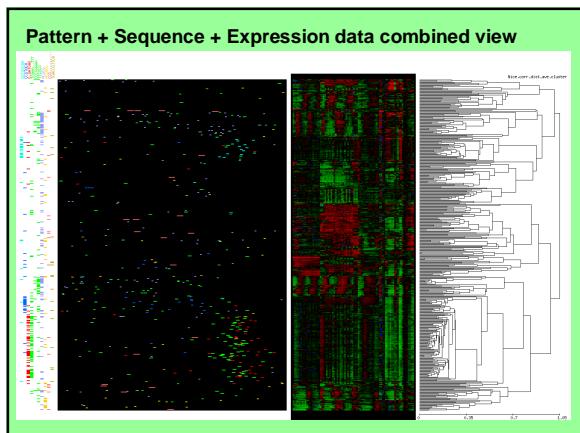
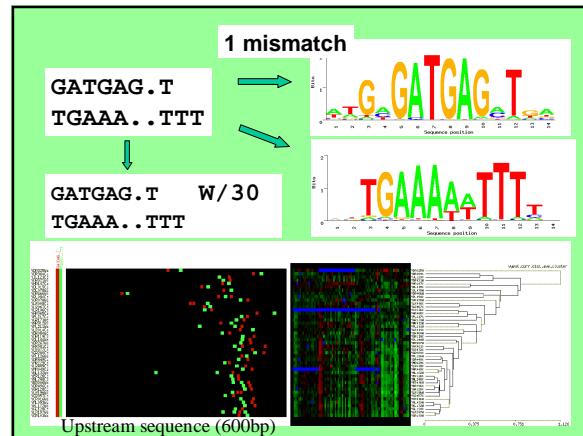
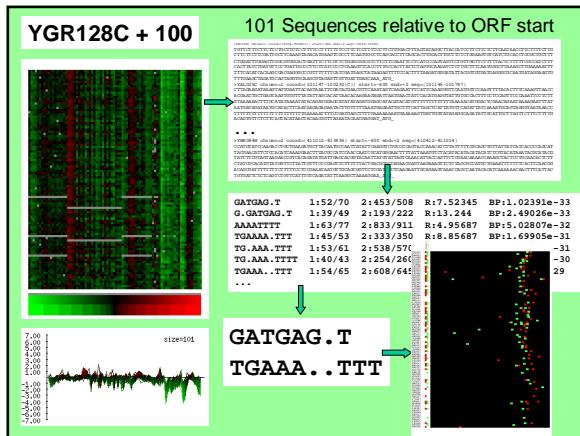


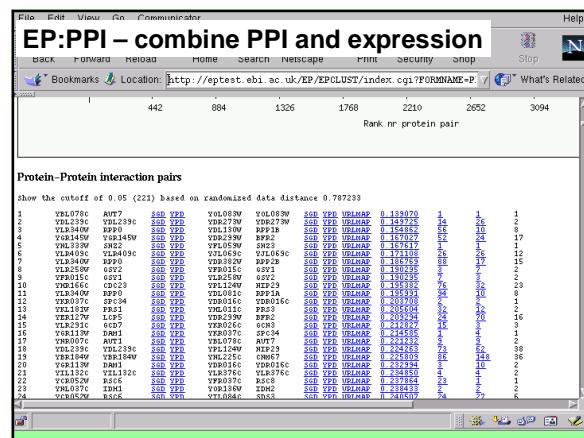
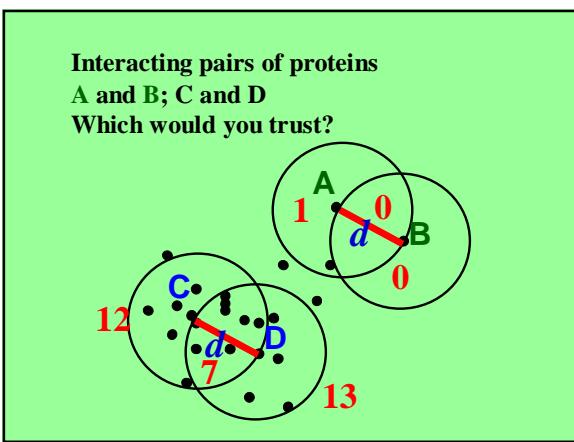
From microarray images to gene expression data







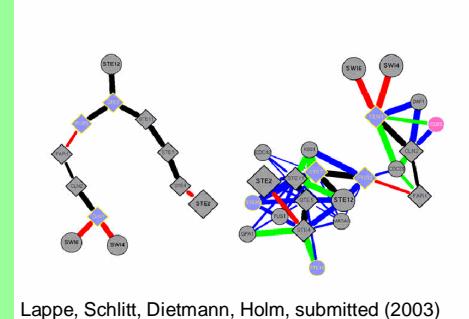




Text mining

- Teadusartiklid sisaldavad tohutul hulgal teavet
- Kuidas seda süsteemiliselt kokku koguda ja esitada?
- Tuvasta Medline abstraktides esinevad geenide nimed ja seosed nende vahel
- Moodusta selle info põhjal graafilised esitlused seoste kohta

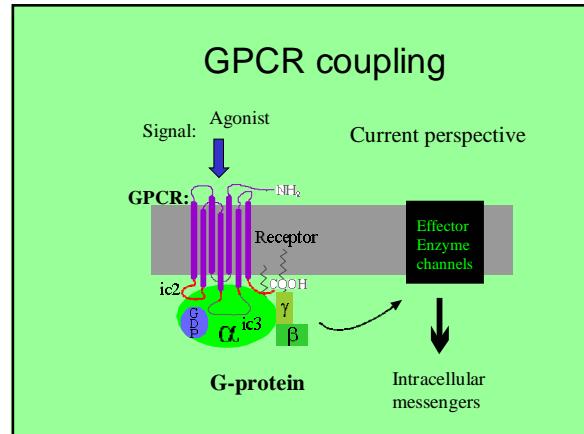
Signal transduction pathway from Text mining and experimental data



Lappe, Schlitt, Dietmann, Holm, submitted (2003)

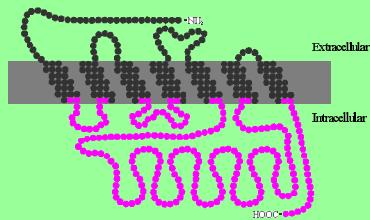
GPCR

- G-valk retseptorid on ühed tähtsamad ravimite sihtmärgid
- Küsimus – kas saab ennustada retseptori põhjal millise signaalülekande raja see käivitab?
- Millised G-valgud seonduvad selle GPCR valguga?

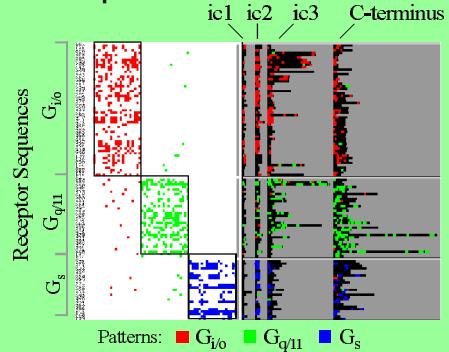


Our Computational Approach

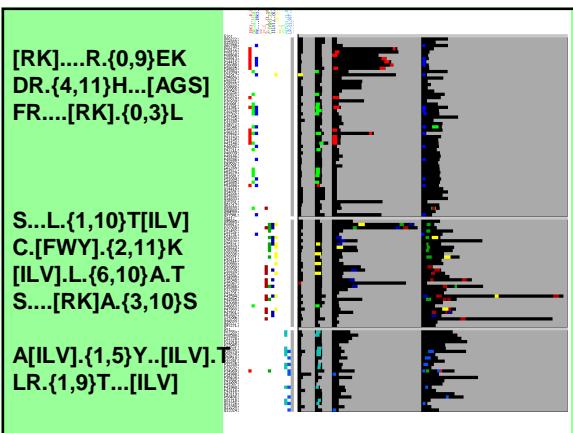
- Using a new membrane topology prediction algorithm (designed specifically for GPCRs), we constrained our pattern search to the intracellular domains of ≈ 100 receptor sequences with well-characterised, and non-promiscuous coupling (split into G_s , $G_{i/o}$ and G_{q11})



Receptor Match Positions

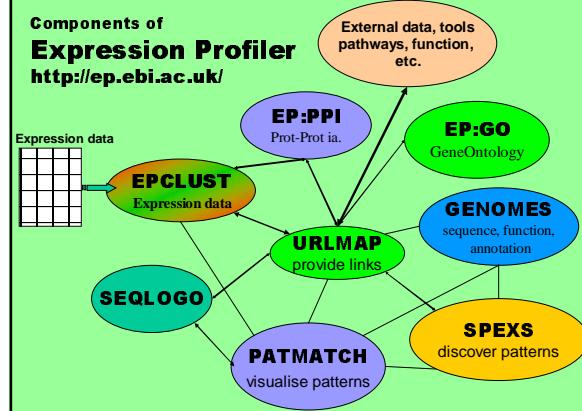


Croning, Vilo, Möller, ISMB 2001



Components of Expression Profiler

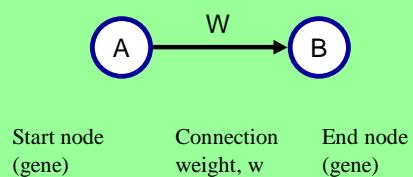
<http://ep.ebi.ac.uk/>



Networks

- Graphical models
- Directed labelled graph
- Nodes \Leftrightarrow genes
- Arcs/Edges \Leftrightarrow relationships
- Labels \Leftrightarrow types of relationships

Graph drawing

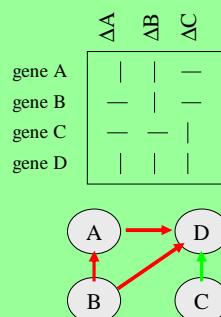


Different interpretation of arcs

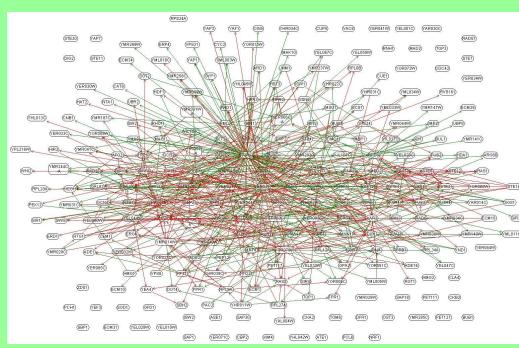
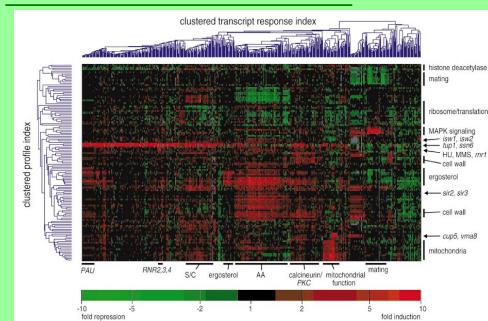
- Edges can have different meanings, hence different networks
- Binding site for A is in front of B
- Proteins A and B interact
- Deletion of gene A affects expression of B (is somewhere in regulation cascade)
- “Literature” mentions genes together

Features/distributions that do not depend on discretisation thresholds

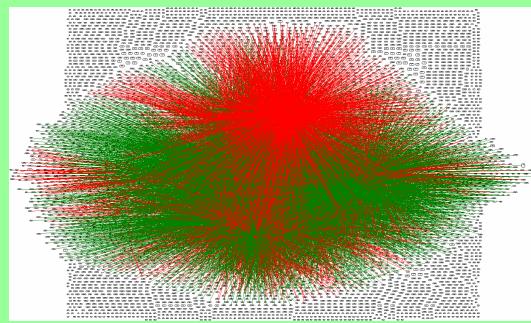
- Visual inspection, biological interpretation
- General statistics and features of the graphs
- Indegree/Outdegree
- Complexity of the networks
- What is the modularity?
 - How many components?
 - Deletion of hot-spots, does it break the net?

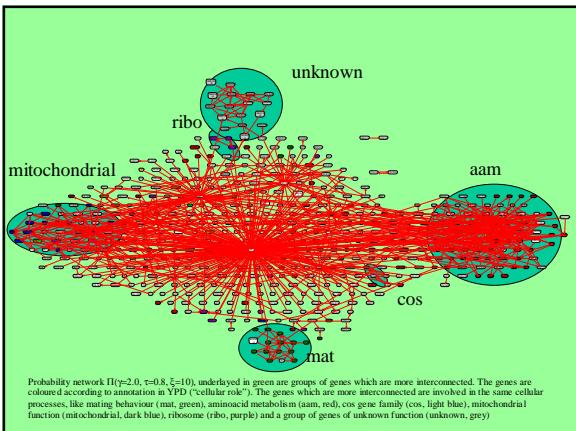
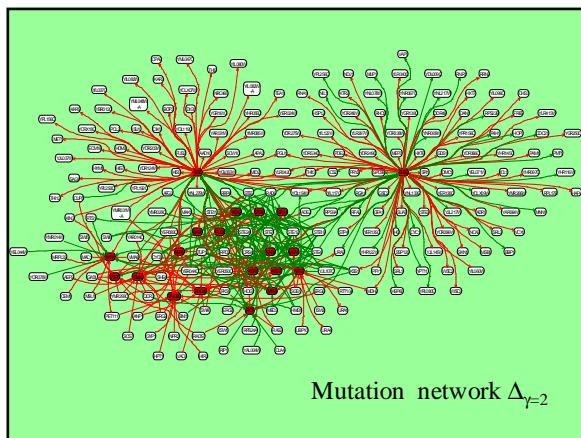
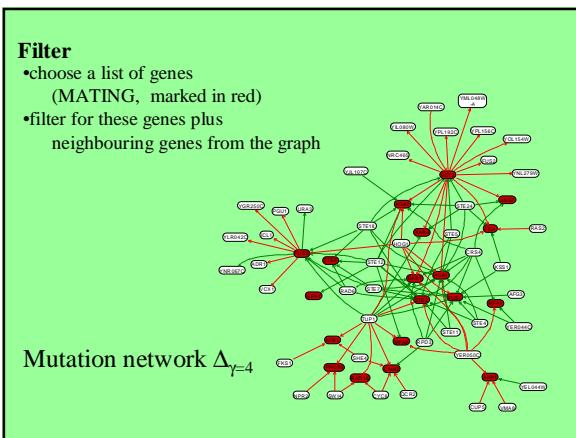


Hughes, T. R. et al: “Functional Discovery via a Compendium of Expression Profiles”, Cell 102 (2000), 109-126.



A complete graph





Kokkuvõttes

- (Molekulaar-) Bioloogilisi andmeid on palju erinevaid tüüpe
- Bioloogiliste andmete mahud kasvavad praegu eksponentsiaalise kiirusega
- Arvutianalüs on ainuke reaalne viis neid andmeid ära kasutada
- See eeldab loomingulist lähenemist nii küsimusepüstituses kui ka tehniliselt heade algoritmide kasutust või väljamõlemist

Bio-andmete kaevandamine

- Millised bio-andmed on olemas ja kuidas neid võiks kasutada
- Andmete kogumine, puastamine, ettevalmistamine ja ühendamine
- Analüs arvutil (peamine algoritm)
- Tulemuste esitamine, reeglite tuvastamine, visualiseerimine
- Tulemuste tõlgendamine

Acknowledgements

Alvis Brazma + the EBI microarray team
Misha Kapushesky, EBI (EP development)
Patrick Kemmeren, Frank Holstege, Utrecht U. (PPI)
Esko Ukkonen, Kimmo Palin, Helsinki Univ.

Meelis Kull Tanel Kaart Hedi Peterson Kristo Käärmann jne.	Maido Remm Reidar Andreson ...
--	---