

Sequence analysis

Esko Ukkonen

University of Helsinki
Helsinki University of Technology



Overview

- n sequences are everywhere
- n combinatorial pattern matching
- n pattern discovery in sequences
- n dynamic programming, automata theory, advanced data structures, probabilistic modeling

Background

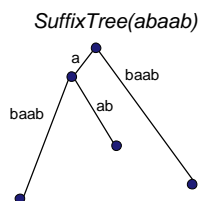
- n algorithms on strings and biological sequence analysis studied by the group since 1980
- n lots of our earlier results appear in textbooks in the field
 - n Gusfield, Cambridge Univ Press 1997
 - n Crochemore and Rytter, Oxford Univ Press 1994
 - n Navarro and Raffinot, Cambridge Univ Press 2002
 - n Smyth, Addison-Wesley 2003
 - n Wikipedia

Main tasks in the theme

- S.1 String algorithms
 - Ukkonen, Mannila, Toivonen
- S.2 Finding orders from data
 - Mannila, Hyvärinen, Ukkonen
- S.3 Sequence segmentation and structure
 - Mannila, Toivonen, Ukkonen
- S.4 Kernel algorithms for sequences
 - Kivinen, Ukkonen

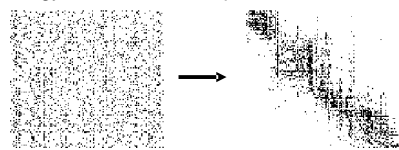
Example result (suffix-array construction)

- n direct construction of a suffix array in linear time
 - n immediately included in teaching materials internationally
 - n J. Kärkkäinen, P. Sanders, S. Burkhardt: *Linear work suffix array construction*, J. ACM (in press), ICALP

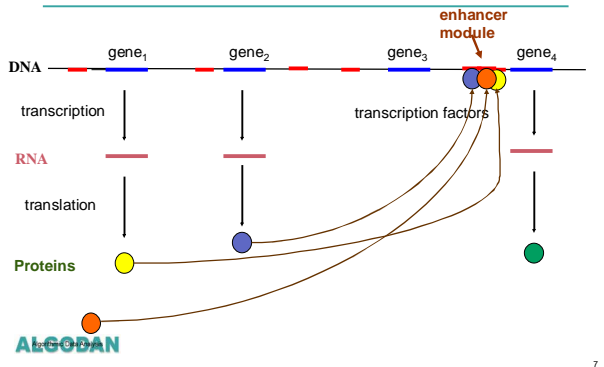


Example (finding orders from data)

- n Seriation problem for 0-1 data in paleontology:
 - n Find an ordering of the rows such that the 1s are as consecutive as possible
- n No errors $\hat{=}$ polynomial
- n Errors $\hat{=}$ NP-hard
- n Spectral techniques; MCMC; finding partial orders (Paleobiology 2006; PLoS Comput Biol 2006; KDD 2006)



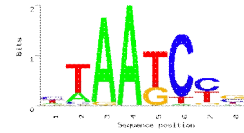
Example (finding gene enhancer modules)



Binding affinity matrices

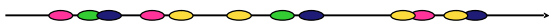
- Transcription factor binding sites represented by affinity matrices
- Discovered:
 - Computationally
 - Traditional wet lab
 - Microarrays

9	11	49	51	0	1	1	4
19	3	0	0	0	45	25	16
5	1	2	0	17	0	4	21
18	36	0	0	34	5	21	10



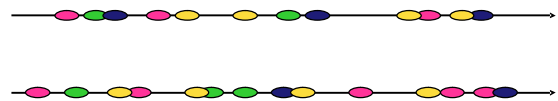
ALGODAN

Good binding sites



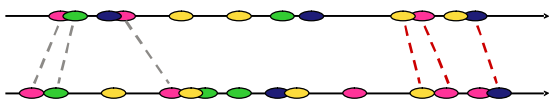
ALGODAN

Clustering and conservation



ALGODAN

Clustering and conservation

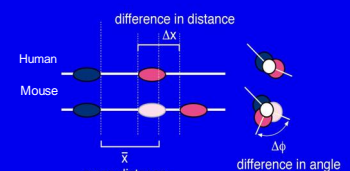


ALGODAN

Computational identification of enhancer elements

- Preserved in evolution:

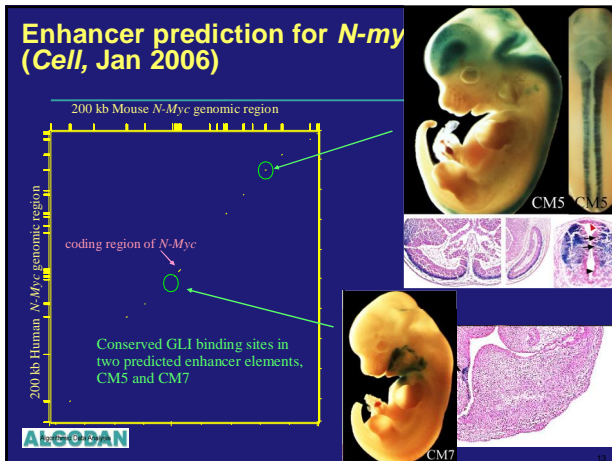
- Affinities of functional cis-elements.
- Spatial arrangement of elements within a module.



$$\text{Score} = \lambda \Delta G_T - \mu \bar{x} - \frac{\sqrt{\Delta X^2 + \Delta \phi^2}}{2\bar{x}}$$

relative weights

affinity clustering conservation



Regulatory modules in gene regulatory networks

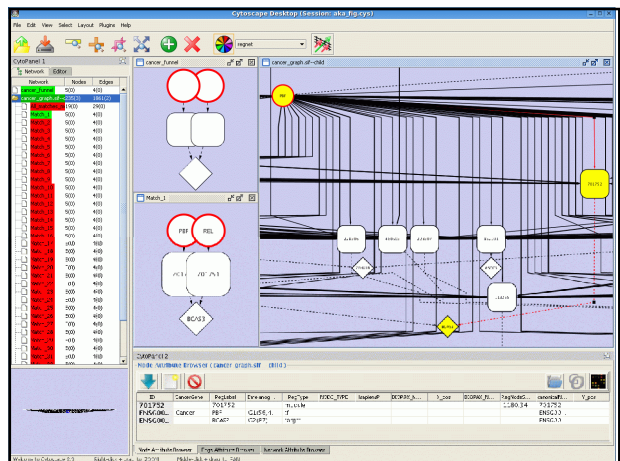
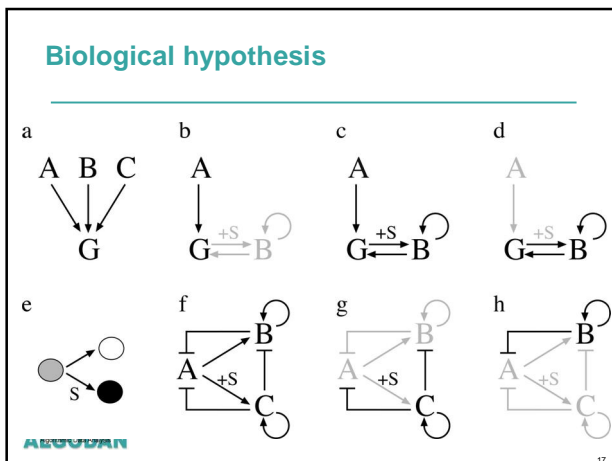
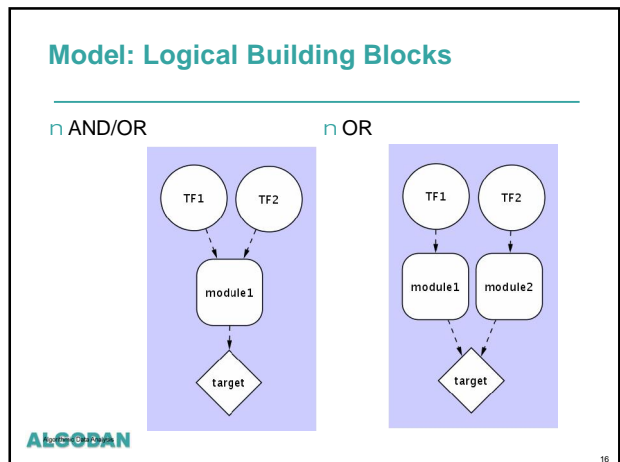
- Still today, organ specific growth control remains poorly understood and difficult to study with traditional genetics
- Now we can start to investigate the "program of development" in top-down fashion instead of bottom-up
 - Soon (within 10 years) the binding specificities of all TFs will be known (?)
 - Tools like EEL can be used to build genome-wide predictions of transcriptional regulation
- Central idea: **Developmental regulatory circuits are embedded in the genome-wide transcriptional network but extracting this knowledge requires new computational tools**

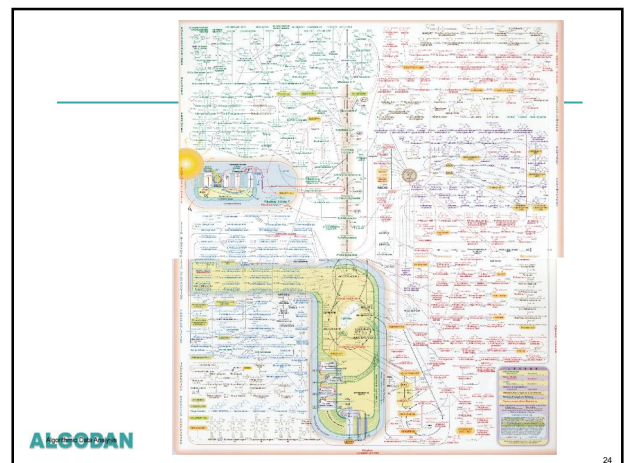
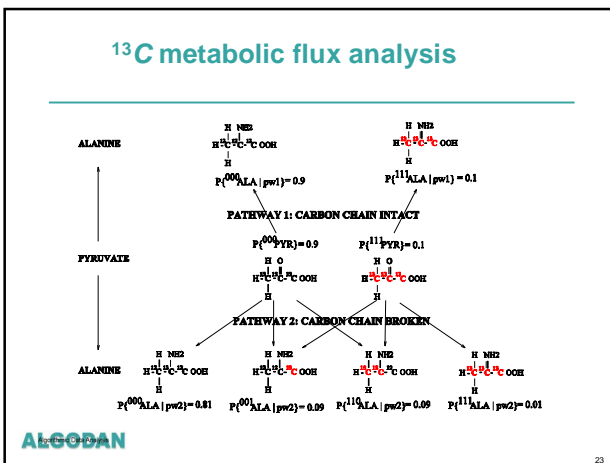
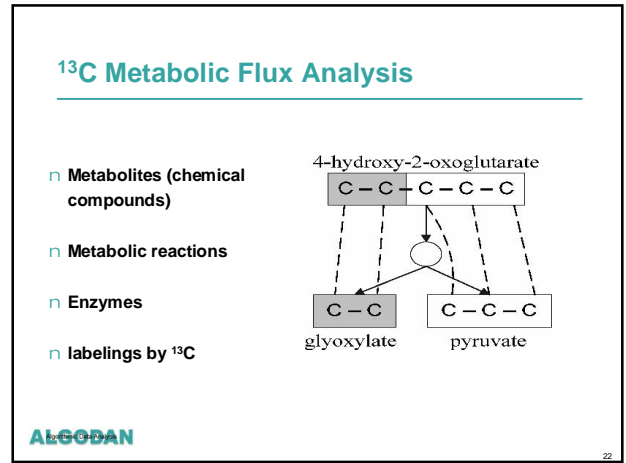
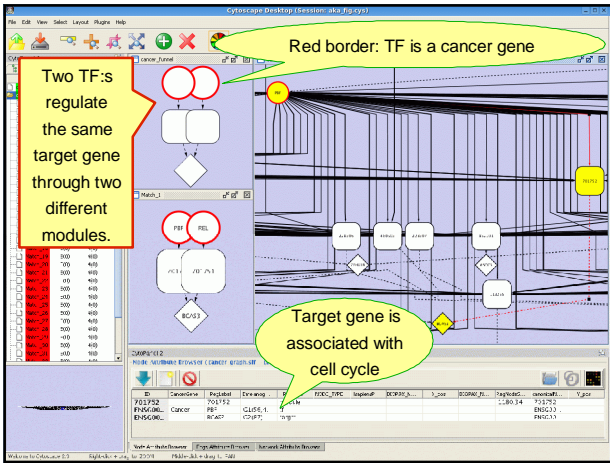
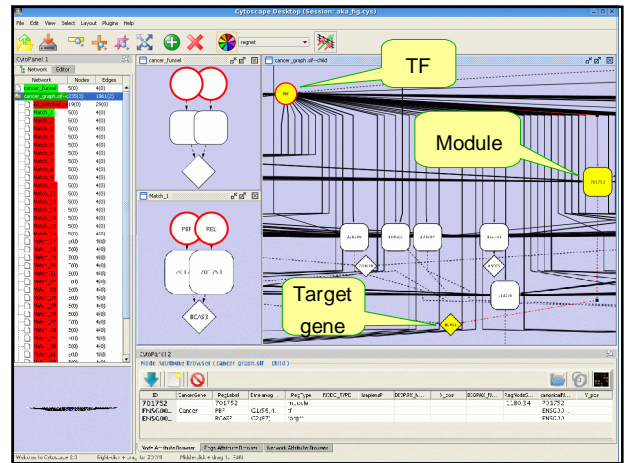
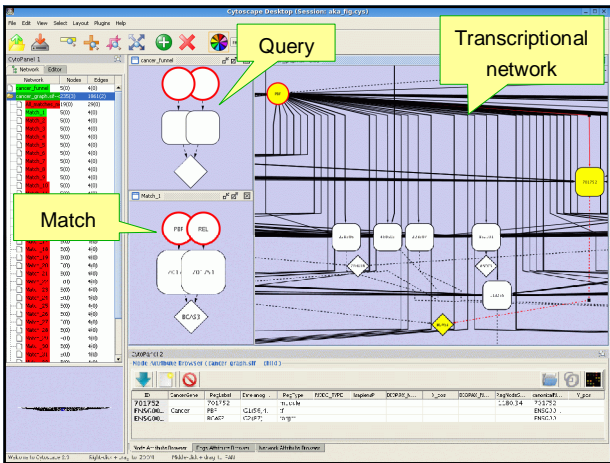
ALGODAN

Model

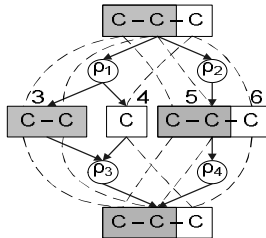
- The control system for growth consists mainly of
 - Transcription factors (TFs)
 - proteins that bind to DNA
 - Regulatory modules
 - clusters of TF binding sites in the regulatory areas of genes that specify the interactions of TFs with the regulated genes.
- A module acts as a **logic gate** which
 - input is the expression levels of the TFs and
 - output is the expression of the regulated gene
- The regulatory interactions form a **complex logic circuit**
 - One TF usually interacts with many modules and the genes encoding TFs are themselves controlled using the same mechanism

ALGODAN





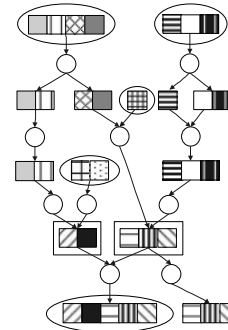
¹³C Metabolic Flux Analysis



ALGODAN

25

¹³C Metabolic Flux Analysis



ALGODAN

26

¹³C Metabolic Flux Analysis

- n Problem: Estimate the fluxes of metabolites
- n Analogous problem: How much gin and how much tonic in gin-tonic?
- n Measure the alcohol contents => solution of the problem
 - n Alcohol content \approx isotopomer content
 - n gin: 40 %, tonic: 0 %, gin-tonic: 10 %
 - $V_g + V_t = V_s$
 - $0,4 \cdot V_g + 0 \cdot V_t = 0,1 \cdot V_s$
 -
 - gin-tonic: $\frac{1}{4}$ gin, $\frac{3}{4}$ tonic



ALGODAN

27

Future goals

- n Indexing sequential data for approximate searches?
- n Distance functions for sequences?
 - n General theoretical framework, efficient evaluation algorithms, complexity bounds, relations between distances
 - n Application specific distances
 - n Generalizations to different generalized sequences (XML, images, music, ...)
- n Finding structure and signals in sequences
 - n Supervised and unsupervised learning of signals
 - n Statistical significance of the findings

ALGODAN

28

Future goals (cont.)

- n How does the program encoded in genomes work?
- n Integrated analysis of genomic, proteomic, and metabolic data
 - n metabolic modeling (new important area in bioinformatics)
 - n metabolic fluxes
- n Challenge problems
 - n Learning network structures
 - n Finding structure and knowledge in natural language data
 - n The vocabulary, grammar and history of genomes

ALGODAN

29