

Advanced Bioinformatics

Biostatistics & Medical Informatics 776

Computer Sciences 776

Spring 2011

Mark Craven

Dept. of Biostatistics & Medical Informatics

Dept. of Computer Sciences

craven@biostat.wisc.edu

www.biostat.wisc.edu/bmi776/

Agenda Today

- course information
- overview of topics
- introductions

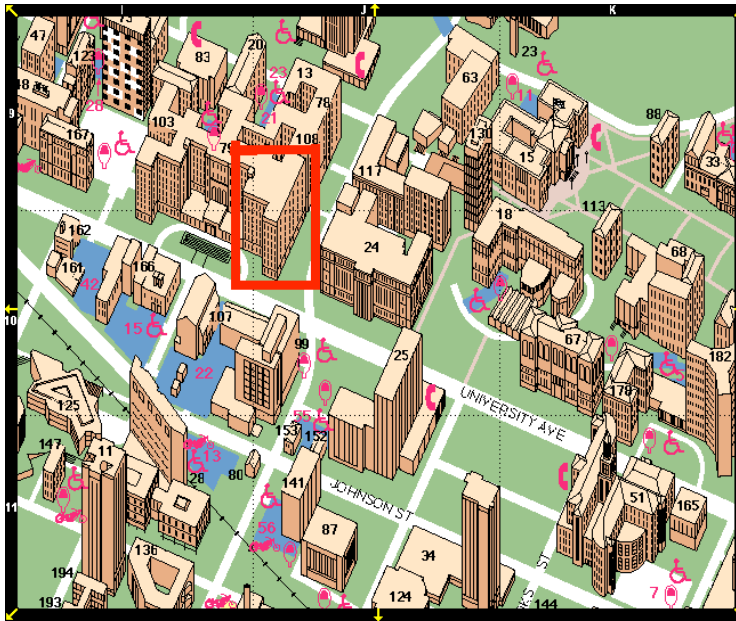
Course Web Site

- www.biostat.wisc.edu/bmi776/
- syllabus
- readings
- tentative schedule
- lecture slides in PDF
- homework
- mailing list archive
- etc.

Your Instructor: Mark Craven

- email:
 craven@biostat.wisc.edu *or*
 craven@cs.wisc.edu
- office hours: TBA
 room 6730, Medical Sciences Center
- my home department is Biostatistics & Medical Informatics, and I have an affiliate appointment in Computer Sciences
- research interests: machine learning, gene regulation and cellular networks, biomedical text mining, probabilistic models, time series

Finding My Office: 6730 Medical Sciences Center



- confusing building
- best bet: enter at door marked *420 North Charter*

Course Requirements

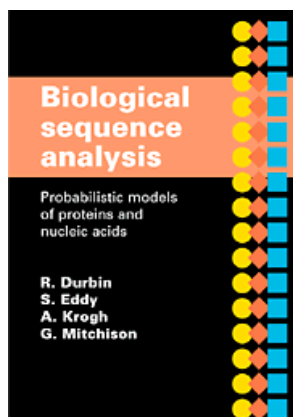
- 4 or so homework assignments: ~20%
 - written exercises
 - programming (in Java, C++, C, Perl, Python) + computational experiments (e.g. measure the effect of varying parameter x in algorithm y)
- 4 or so paper critiques: ~20%
 - major strength of approach
 - major weakness
 - what would you do next
- project: ~25%
- final exam: ~ 25%
- class participation: ~10%

Participation

- Take advantage of the small class size!
- do the assigned readings
- show up to class
- don't be afraid to ask questions

Course Readings

- mostly articles from the primary literature (scientific journals, etc.)
- must be using a UW IP address to download some of the articles
- *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*. R. Durbin, S. Eddy, A. Krogh, and G. Mitchison. Cambridge University Press, 1998.



Computing Resources for the Class

- Linux workstations in Dept. of Biostatistics & Medical Informatics
 - no “lab”, must log in remotely
 - most of you have accounts?
 - two machines
 - mi1.biostat.wisc.edu
 - mi2.biostat.wisc.edu
- CS department usually offers UNIX orientation sessions at beginning of semester
- the “CS 1000” UNIX tutorial
 - online at <http://www.cs.wisc.edu/csl/cs1000/>

The Class Mailing List

- bmi776-1-s11@lists.wisc.edu
- you will be automatically subscribed
- check your mail daily or have it forwarded to an account where you do

Major Topics to be Covered (the task perspective)

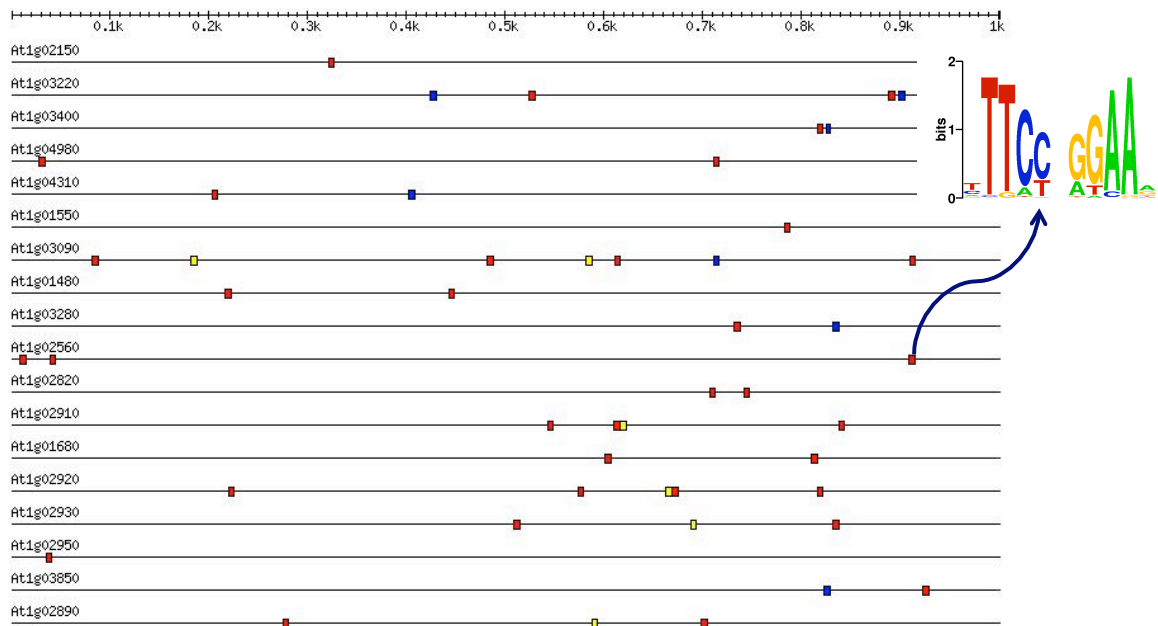
- modeling of motifs and *cis*-regulatory modules
- gene finding
- large-scale and whole-genome sequence alignment
- RNA sequence and structure modeling
- modeling cellular networks
- protein structure prediction
- biomedical text mining
- genotype analysis and association studies

Major Topics to be Covered (the algorithms perspective)

- Gibbs sampling and EM
- HMM structure search
- duration modeling and semi-Markov models
- pairwise HMMs
- interpolated Markov models and back-off methods
- parametric alignment
- tries and suffix trees
- sparse dynamic programming
- Markov random fields
- stochastic context free grammars
- Bayesian networks and module networks
- active learning
- branch and bound search
- conditional random fields
- etc.

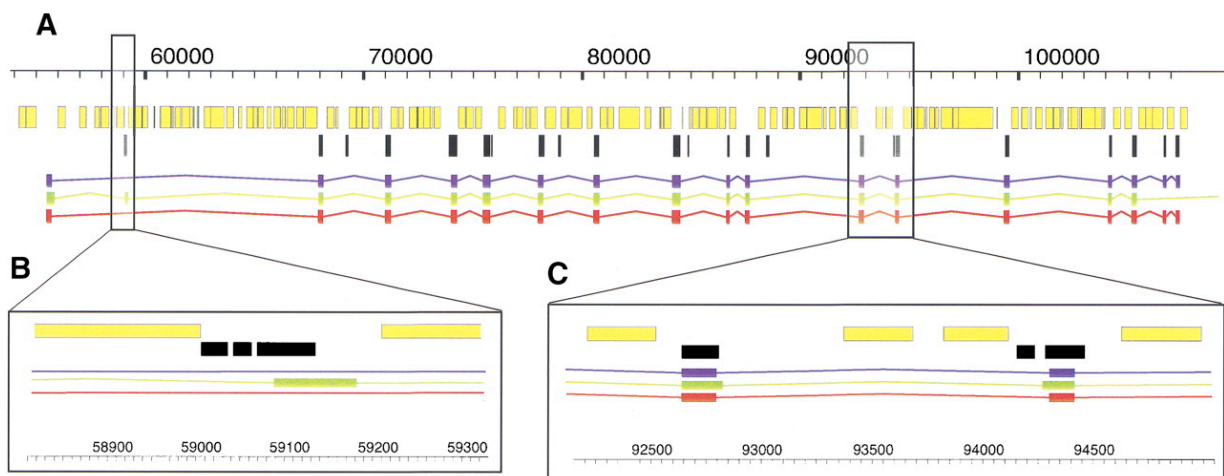
Motif and CRM Modeling

What sequence motifs do these promoter regions have in common?



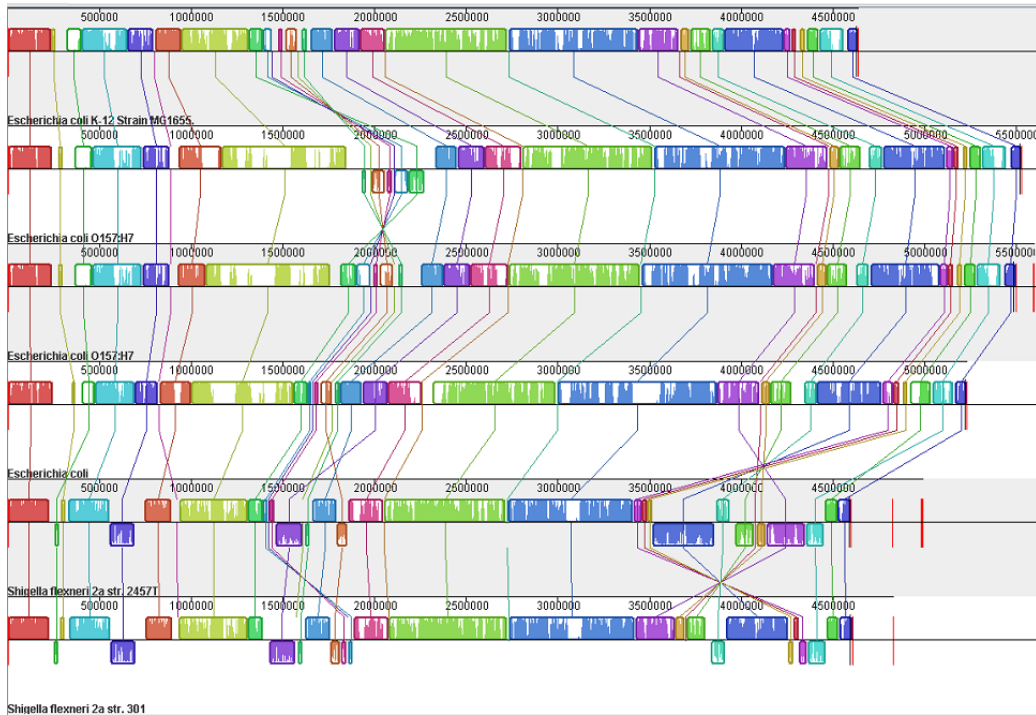
Gene Finding

Where are the genes in this genome, and what is the structure of each gene?



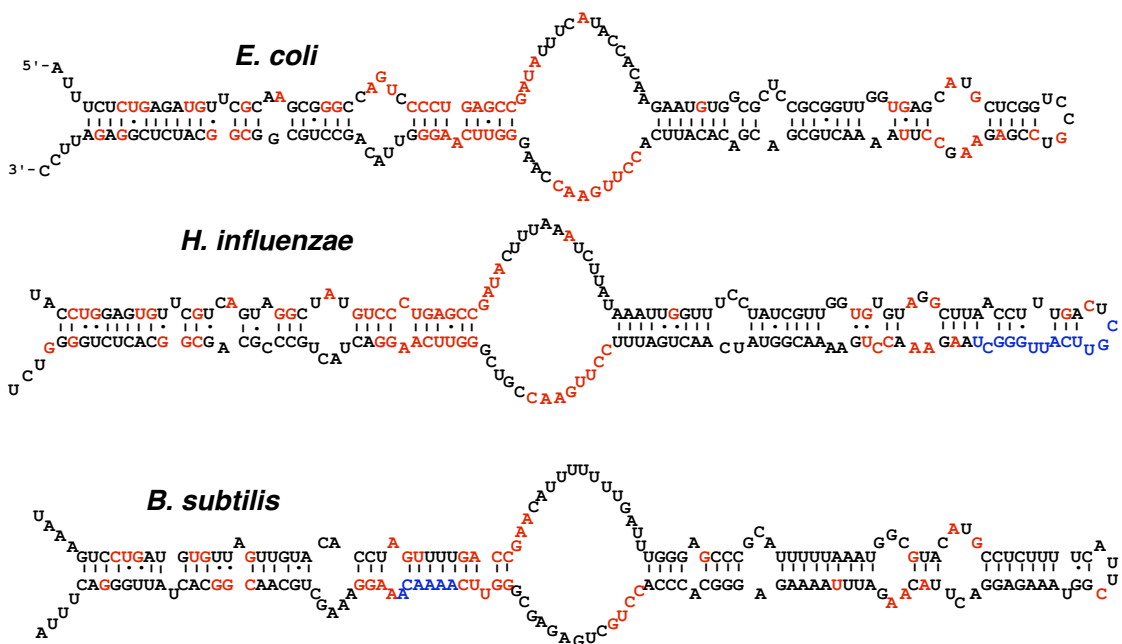
Large Scale Sequence Alignment

What is the best alignment of these 5 genomes?



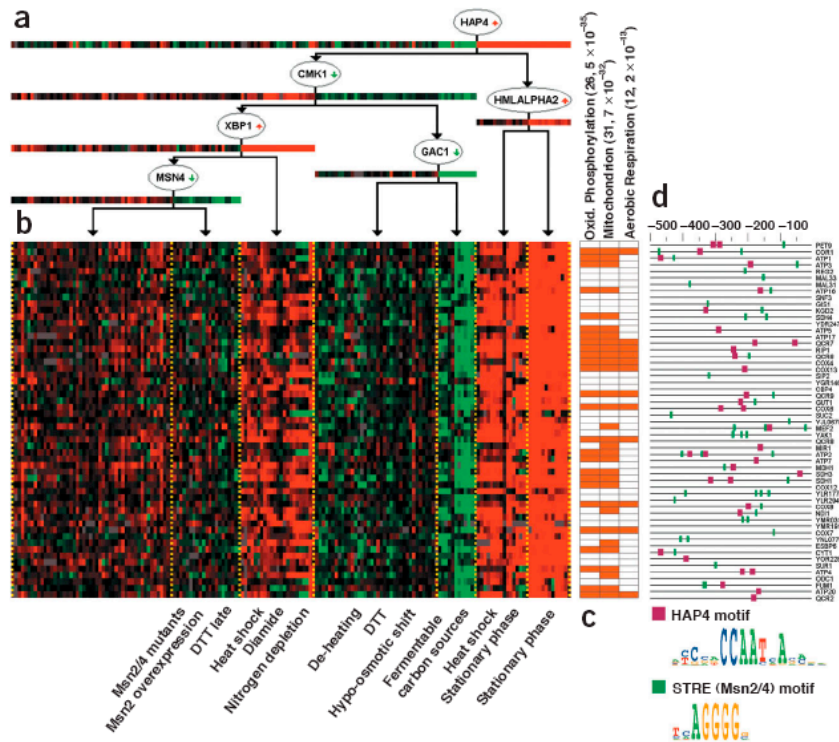
RNA Sequence and Structure Modeling

Given a genome, how can we identify sequences that encode this RNA structure?



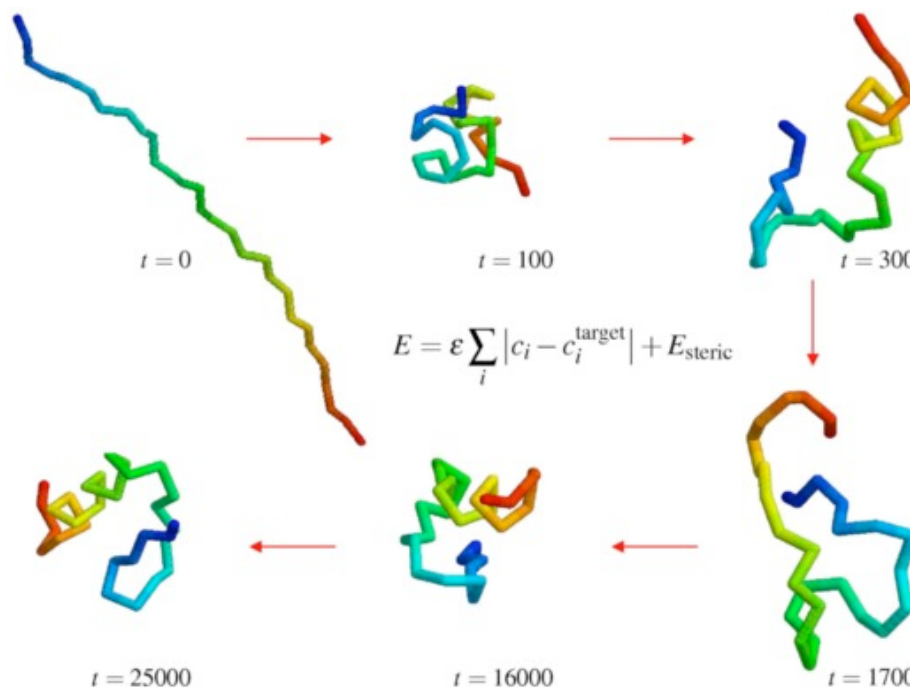
Inferring Intracellular Networks

Can we automatically infer models of regulatory networks from data?



Protein Structure Prediction

Can we predict the 3D shape of a protein from its sequence?



Biomedical Text Mining

Gene Detail

Symbol: **Fut4**
Name: **fucosyltransferase 4**
ID: MGI:95594

Synonyms: 3-fucosyl-N-acetyl-lactosamine, 3-fucosyl-N-acetyl-lactosamine, alpha(1,3) fucosyltransferase, myeloid specific, FAL, FuT-IV, SSEA-1

Map position: Chromosome 9, 3.0 cM

Mammalian orthology: human: rat (Mammalian Orthology)

Sequences: human: rat (Mammalian Orthology)

Phenotypes: All phenotypic alleles(1): Targeted(1)

Polymorphisms: RFLP(1)

Gene Ontology (GO) classifications: Process: [protein amino acid glycosylation](#), [Golgi apparatus, integral to membrane...](#); Component: [Golgi apparatus, integral to membrane...](#); Function: [fucosyltransferase activity](#), [transferase activity...](#)

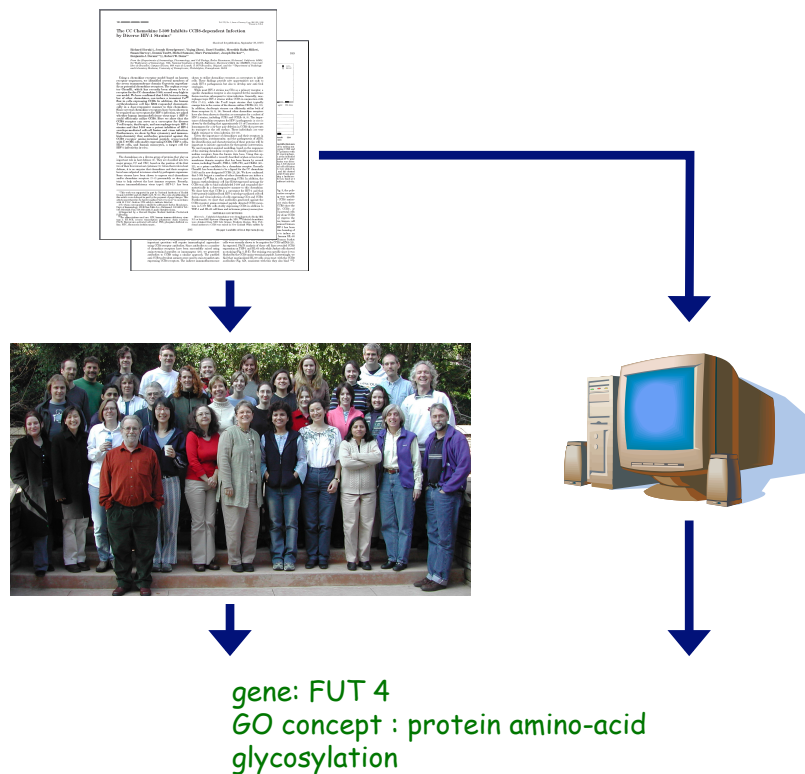
Expression: Theiler Stage 1, 2, 3, 5, 9, 11, 13, 15, 17, 19, 20, 21, 22, 23, 24, 28; Tissues(61); Assay Type: Immunohistochemistry (70), Assays(4); GXD literature index(20); cDNA source data(2)

Other database links: DoTS: DT:40171675, DT:91334210; UniGene: 63450; ENSEMBL: ENSMUSG00000049307; LocusLink: 14345; NIA Mouse Gene Index: NAP015586-001; Entrez Gene: 14345

Protein domains: InterPro ID Description: IPR001503 Glycosyl transferase, family 10; Graphical View of Protein Domain Structure

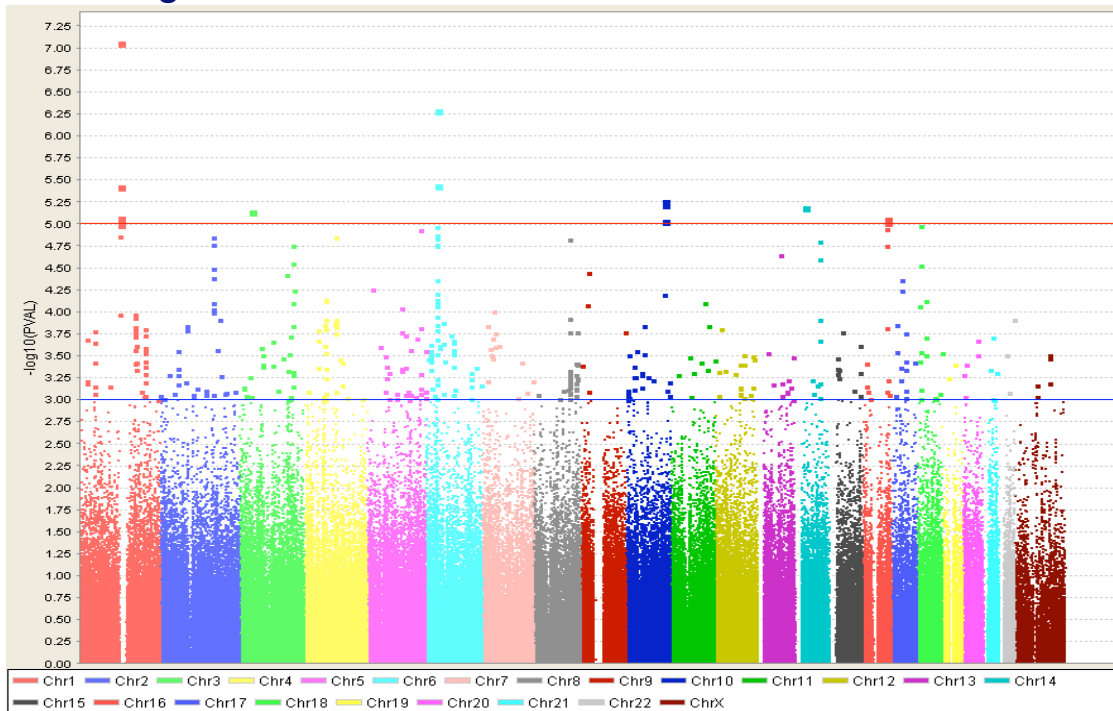
Can we partially automate the process of curating genomic databases?

Biomedical Text Mining



Genome-wide Association Studies

Which genes are involved in diabetes?



Type 2 diabetes association P values by chromosome (386,731 markers). The x-axis is the genomic position by chromosome 1-22 and X (by color), and the y-axis is the negative base 10 logarithm of the P value.

Reading Assignment

- Bailey and Elkan, *ISMB '95*
- Lawrence et al., *Science '93*
- available on the course web site