# Advanced Bioinformatics
## Biostatistics & Medical Informatics 776
## Computer Sciences 776
## Spring 2012

Colin Dewey

Dept. of Biostatistics & Medical Informatics

Dept. of Computer Sciences

cdewey@biostat.wisc.edu

www.biostat.wisc.edu/bmi776/

# Agenda Today

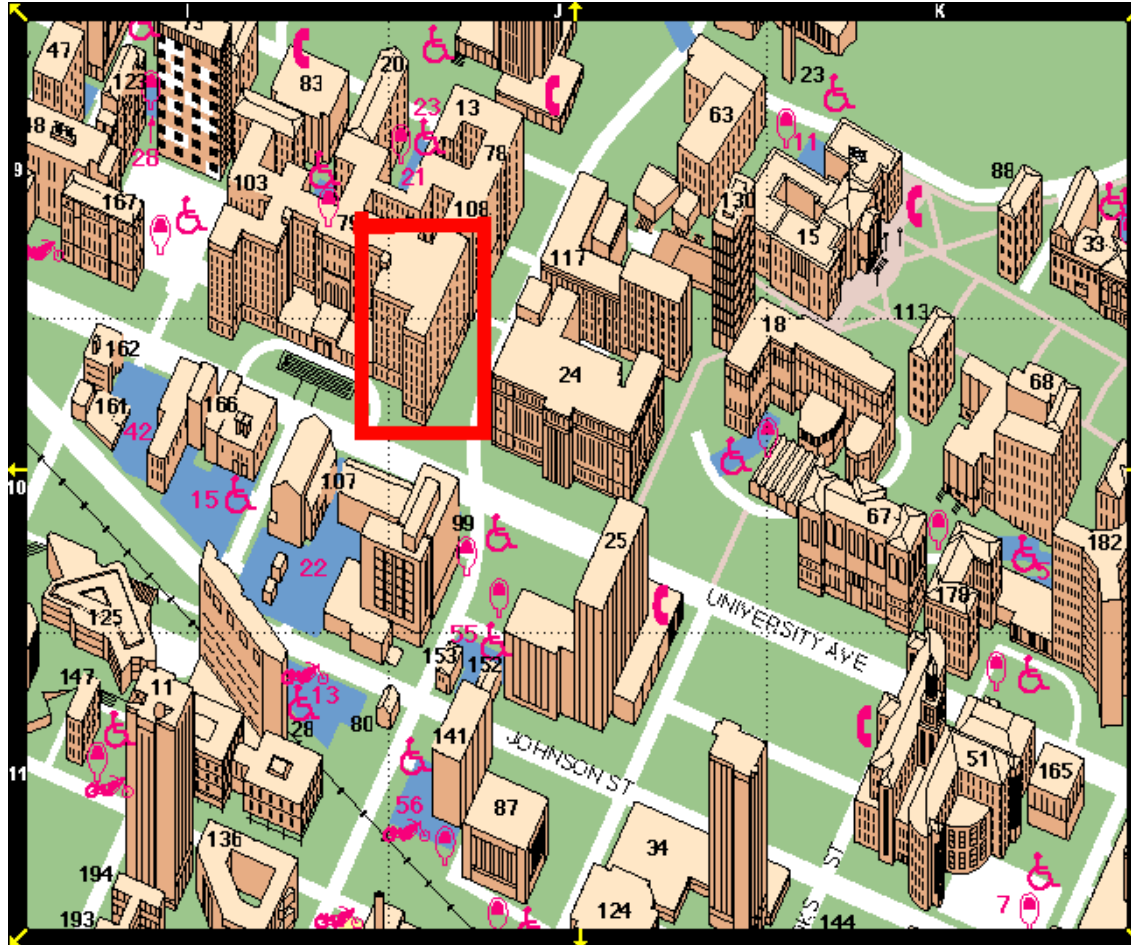- course information
- overview of topics
- introductions

# Course Web Site

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

# Your Instructor: Colin Dewey

- email:
  cdewey@biostat.wisc.edu

- office hours: Wed 9:30-10:30am, Thu 11:00am-12:00pm
  room 5785, Medical Sciences Center

- my home department is Biostatistics & Medical Informatics, and I have an affiliate appointment in Computer Sciences

- research interests: probabilistic modeling, biological sequence evolution, analysis of "next-generation" sequencing data (RNA-Seq in particular)

# Finding My Office:
# 5785 Medical Sciences Center



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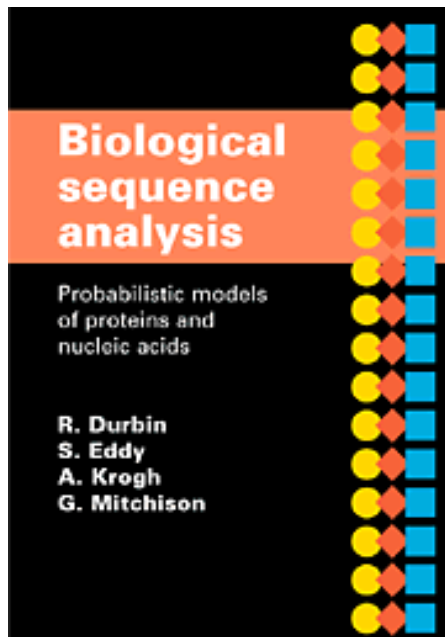
# Course Requirements

- 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 (can use WiscVPN "On Campus" profile)

- *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-s12@lists.wisc.edu
- you will be automatically subscribed
- check your mail daily or have it forwarded to an account where you do
  - mailing list has your @wisc.edu address

# What you should get out of this course

- An understanding of the major problems in computational molecular biology

- Familiarity with the algorithms and statistical techniques for addressing these problems

- At the end you should be able to:
  - Read the bioinformatics literature
  - Apply the methods you have learned to other problems both within and outside of bioinformatics

# Major Topics to be Covered
## (the task perspective)

- modeling of motifs and *cis*-regulatory modules
- identification of transcription factor binding sites
- gene finding
- transcriptome quantification and assembly
- RNA sequence and structure modeling
- modeling biological sequence evolution
- large-scale and whole-genome sequence alignment
- modeling the evolution of 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
- branch and bound search
- conditional random fields
- etc.

# Motif and CRM Modeling

What sequence motifs do these promoter regions have in common?

# Experimental binding site prediction with second-generation sequencing data (ChIP-Seq)



(Park, Nat Rev Genet, 2009)

# Gene Finding

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

# Transcriptome analysis with RNA-Seq

# Modeling biological sequence evolution

```
Cow        …ATC-AT…
Dog        …AGC-AT…
Mouse      …ACC-AT…
Rat        …ACCGAT…
Macaque    …A-CGTT…
Chimp      …A--GTT…
Human      …A-CGTT…
Ancestor …??????…
```



…??????…

…ATCCAT…    …AGCGAT…    …ACCGAT…    …ACCGAT…    …ATCGTT…    …ATCGTT…    …ATCGAT…

# 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?

# Modeling cellular network evolution



g  RNA metabolism

yeast          worm          fly

# Protein Structure Prediction

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



$$E = \varepsilon \sum_i \left| c_i - c_i^{\text{target}} \right| + E_{\text{steric}}$$

# Biomedical Text Mining



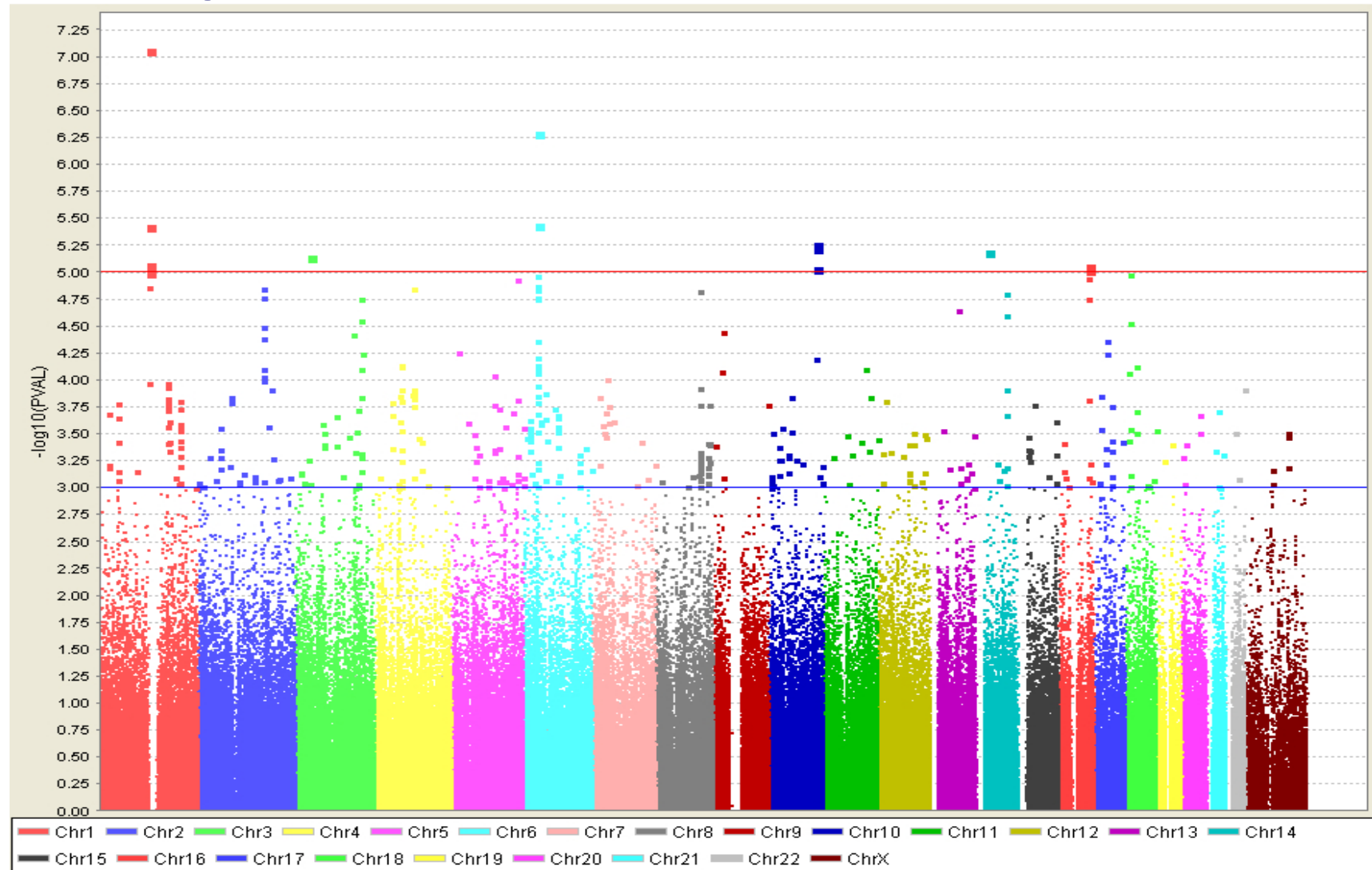Can we partially automate the process of curating genomic databases?

# Biomedical Text Mining

gene: FUT 4
GO concept : protein amino-acid glycosylation

# 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