

CH391L Bioinformatics (course # 52990)

Spring Semester, 2011

Instructor: Edward Marcotte marcotte@icmb.utexas.edu
Office: MBB 3.210AA Phone: 471-5435
Course lectures: **T/Th 12:30 – 2:00 PM WEL 3.402**
Office hours: Wednesdays 2:00 – 3:00 PM MBB 3.210AA
TA: Taejoon Kwon taejoon.kwon@mail.utexas.edu
TA Office hours: Tuesday/Thursday 10:00 – 11:00 AM MBB 3.210A TA Phone: 232-3919
Course web page: **<http://www.marcottelab.org/index.php/CH391L>**

Open to graduate students and upper division undergraduates in natural sciences and engineering.

Prerequisites: Basic familiarity with molecular biology, statistics & computing, but realistically, it is expected that students will have extremely varied backgrounds.

An introduction to computational biology and bioinformatics. The course covers typical data, data analysis, and algorithms encountered in computational biology. Topics will include introductory probability and statistics, basics of programming, protein and nucleic acid sequence analysis, genome sequencing and assembly, protein structure prediction, analysis of DNA microarray data, data clustering, biological pattern recognition, and biological networks.

** Note that this is not a course on practical sequence analysis or using web-based tools. Although we will use a number of these to help illustrate points, the focus of the course will be on learning the underlying algorithms and exploratory data analyses and their applications. **

Most of the lectures will be from research articles and handouts, with some material from the...

Recommended text *Biological sequence analysis*
(for sequence analysis): R. Durbin, S. Eddy, A. Krogh, G. Mitchison
Cambridge University Press
Avail. from Amazon.com (\$51.30, 5-9 days delivery)

For non-molecular biologists, I highly recommend (really!) *The Cartoon Guide to Genetics* (Gonick/Wheelis)
For biologists rusty on their stats, *The Cartoon Guide to Statistics* (Gonick/Smith) is also very good.

Some online references:

Online bioinformatics course: <http://lectures.molgen.mpg.de/>
Bioinformatics algorithms: <http://sapc34.rdg.ac.uk/~andrew/algorithms/>
Bioinformatics resources on the web: <http://zlab.bu.edu/zlab/links.shtml>
Online probability texts: <http://omega.albany.edu:8008/JaynesBook.html>
 <http://www-users.york.ac.uk/~mb55/pubs/pbstnote.htm>
 http://www.dartmouth.edu/~chance/teaching_aids/books_articles/probability_book/pdf.html

No exams will be given. Grades will be based on 4 problem sets (given every 2 weeks and counting 15% each towards the final grade) **and a course project** (40% of final grade), which can be individual or collaborative. If collaborative, cross-discipline collaborations are encouraged. The course project will consist of a research paper or project on a bioinformatics topic chosen by the student (with approval by the instructor) containing an element of independent computational biology research (e.g. calculation, programming, database analysis, etc.). This will be turned in as a link to a web page.

The final project is due on May 3, 2011.

Outline for CH391L **Bioinformatics**

- 1/18 Introduction + next-generation genome sequencing sample collection
- 1/20 A biology primer for non-molecular biologists (DNA, RNA & proteins)
- 1/25 A Perl programming primer for non-programmers
- 1/27 Start biological sequence analysis
- 2/1 NO CLASS
- 2/3 Biological sequence analysis
- 2/1, 3/10 NO CLASS (also no class on Spring Break, 3/15, 3/17)

We'll cover the following topics, in order:

BIOLOGICAL SEQUENCE ANALYSIS

Substitution matrices (BLOSSUM, PAM) & sequence alignment
Protein and nucleic acid sequence alignments, dynamic programming
Sequence profiles
BLAST! (the algorithm)
Biological databases
Markov processes and Hidden Markov Models
Gene finding algorithms

GENOMES, DNA MICROARRAYS, & "BIG BIOLOGY"

Gene finding algorithms & GASP
An introduction to genome sequences & shotgun sequencing
Genome assembly & how the human genome was sequenced
Next- (& next-next-) generation DNA sequencing & the revolution in genomics
An introduction to DNA microarrays and large gene expression data sets
Clustering algorithms, hierarchical, k-means, self-organizing maps, force-directed maps
Classifiers, k-nearest neighbors, Mahalanobis distance
Promoter and motif finding, Gibbs sampling
Principal component analysis and data transformations

BIOLOGICAL NETWORKS & SYNTHETIC BIOLOGY

Biological networks: metabolic, signaling, graphs, regulatory
Properties of biological networks
Network alignment and comparisons, network organization
Analogies between biological networks and electrical circuits
Designing, simulating, and building gene circuits

PROTEIN FOLDING AND STRUCTURE PREDICTION (Optional—we'll vote on this, or consider alternate topics)

Fold recognition, 3D-1D profiles, threading, *ab initio* design

***** FINAL PROJECT DUE on May 3, 2011 *****