

Instructor: Prof. Edward Marcotte marcotte@icmb.utexas.edu Office hours: Wed 11 – 12 MBB 3. 148BA TA: Brendan Floyd bmfloyd@utexas.edu Office hours: Mon 1 – 2/Fri 1:30 – 2:30 NHB 3.400B atrium (or MBB 3.128B) Phone: 512-232-3919

Probably the most important slide today! Course web page: http://www.marcottelab.org/ index.php/BCH394P_BCH364C_2020 This is a graduate student class! It is open to a small # (<10) of upper division undergrads in natural sciences and engineering. UG prerequisites: Biochemistry 339F with a grade of at least B; Computer Science 303E and Statistics and Data Sciences 328M (or Statistics and Scientific Computation 318M, 328M) with a grade of at least C-; and consent of the instructor.

3

An introduction to systems biology and bioinformatics, emphasizing quantitative analysis of high-throughput biological data, and covering typical data, data analysis, and computer algorithms.

Topics will include introductory probability and statistics, basics of Python programming, protein and nucleic acid sequence analysis, genome sequencing and assembly, proteomics, synthetic biology, analysis of large-scale gene expression data, data clustering, biological pattern recognition, and gene and protein networks.

** NOT a course on practical sequence analysis or using web-based tools (although we'll use those too), but rather on algorithms, exploratory data analyses and their applications in high-throughput biology. **

Books

Most of the lectures will be from research articles and slides. For sequence analysis, there will be an **Optional text:**

Biological sequence analysis, Durbin, Eddy, Krogh, Mitchison, Cambridge Univ. Press (available from Amazon, used & ebook)

For biologists rusty on their stats, *The Cartoon Guide to Statistics* (Gonick/Smith) is very good (really!).

We will also be learning some Python programming. I <u>highly</u> recommend...

Python programming for beginners: https://www.codecademy.com/learn/learn-python

5









R	SALIND About - Problems - Statistics - Glossary Search		My Classes - edward.marcotte Log out			
Pro	blems		Bioir	formatics Stri	onghold •	List Tree
Rosalind	is a platform for learning bioinformatics and programming through problem solving. Take a tour to get the ha	ng of how Rosalind works.				
ast win: h	lydratedlizard vs. "Constructing a De Bruijn Graph", 7 minutes ago	P Colored Day Core	roblems: 285	total), users: 60	64, attempts: 101	4640, correct: 5674
DNA	Counting DNA Hustantides	25050	lect Ratio	Questions	Jointons	Explanation
RNA	Transcribing DNA into RNA	31498				
REVC	Complemention a Strand of DNA	28531				
FIB	Rabbits and Recurrence Relations	16249				
GC	Computing GC Content	16729				
HAMM	Counting Point Mutations	18923				
IPRB	Mendel's First Law	10808				
PROT	Translating RNA into Protein	14743				
SUBS	Finding a Motif in DNA	15115				
CONS	Consensus and Profile	8423				
FIBD	Mortal Fibonacci Rabbits	7001				
GRPH	Overlap Graphs	6963				
IEV	Calculating Expected Offspring	6357				
LCSM	Finding a Shared Motif	5909				
LIA	Independent Alleles	3367				
MPRT	Finding a Protein Motif	3708				
MRNA	Inferring mRNA from Protein	5698				
ORF	Open Reading Frames	4399				
PERM	Enumerating Gene Orders	7860				
PRTM	Calculating Protein Mass	7255				
REVP	Locating Restriction Sites	4694				
SPLC	RNA Splicing	5193				
LEXF	Enumerating k-mers Lexicographically	4383				
LGIS	Longest Increasing Subsequence	1924				
LONG	Genome Assembly as Shortest Superstring	2195				
PMCH	Perfect Matchings and RNA Secondary Structures	2062				
		2004				





- By submitting *as your own work* any unattributed material that you obtained from other sources, you have committed plagiarism.
- Copying homework solutions from other students or internet sources is cheating, collusion, and/or plagiarism.
- Software and computer code are legally considered in the same framework as other written works. Copying code directly without attribution is plagiarism.

See the university's official policy on plagiarism here: https://catalog.utexas.edu/general-information/appendices/appendix-c/student-discipline-and-conduct/

- You can use the internet to get *ideas*, programming *suggestions* and *syntax*, but <u>downloading completed answers to</u> <u>assigned questions and submitting these as</u> <u>your own work is cheating/plagiarism</u>.
- <u>Copying entire programs</u> verbatim from marked repositories offering Rosalind homework solutions <u>is cheating and</u> <u>plagiarism</u>.

Similarly, downloading or otherwise obtaining solutions to homework problems from previous students (or Coursehero/similar sites) and turning these in as your own work is cheating, collusion, and/or plagiarism.



Why are we here? (practically, not existentially)



Our current-ish knowledge of human metabolism							
Total number of reactions Total number of metabolites Number of unique metabolites Number of metabolites in extracellular space Number of metabolites in cytoplasm Number of metabolites in mitochondrion Number of metabolites in nucleus Number of metabolites in endoplasmic reticulum Number of metabolites in peroxisome Number of metabolites in lysosome Number of metabolites in Golgi apparatus Number of transcripts Number of unique genes	7,440 5,063 2,626 642 1,878 754 165 570 435 302 317 2,194 1,789						
Nat Biotechnol. 2013 May;31(5):419-25 Updated in Metabolomics 2016 12:109							









Specifically...

We'll cover the following topics, approximately in this order:

BASICS OF PROGRAMMING

Introduction to Rosalind A Python programming primer for non-programmers Rosalind help & programming Q/A

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



Plus, expert guest lectures on:

NGS best practices Overview of mass spectrometry shotgun proteomics Protein 3D structural modeling

Plus, plus: we'll attempt a live demo in-class of nanopore sequencing....

THE FINAL COURSE PROJECT IS DUE by midnight, April 27, 2020

The last 3 class days will be devoted to presenting your projects to the rest of the class.