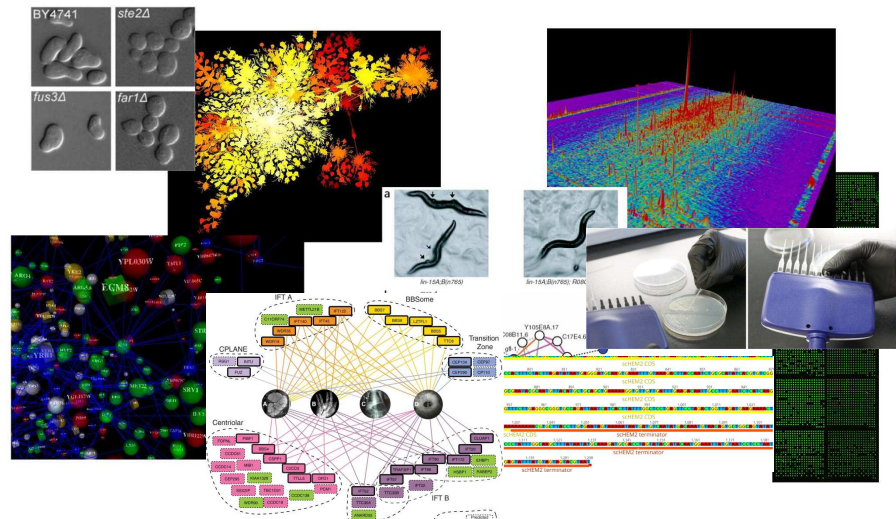


BCH394P/BCH364C Systems Biology & Bioinformatics
(course # 55680/55580)

Spring 2021 Tues/Thurs 11 – 12:30 PM Virtual



1

Instructor: Prof. Edward Marcotte
Zoom office hours: Fri 11 – 12

marcotte@utexas.edu

TA: Vy Dang
Zoom office hours: Mon 3 – 4/Wed 1 – 2

vyqtdang@utexas.edu

Class Slack channel: bch394pbch364c2021.slack.com

The class zoom channel will be posted on Canvas.
It will be the same zoom for class and office hours.

2

Probably the most important slide today!

Course web page:

**[http://www.marcottelab.org/
index.php/BCH394P_BCH364C_2021](http://www.marcottelab.org/index.php/BCH394P_BCH364C_2021)**

This is a graduate student class!

It is open to a small # 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. ****

4

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 highly recommend...

Python programming for biologists:

<https://pythonforbiologists.com/introduction/>

5

Grading

No exams. Instead, grades will be based on:

- **Online programming homework**
(10 points each and counting 30% of the final grade)
- **3 problem sets**
(15 points each and counting 45% of the final grade)
- **A course project** that you will develop over the semester & present in the last 2.5 days of class (25% of final grade)

The course project will consist of a research 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.) turned in as a web URL (20%) and presented in class (5%).

The project will be emailed as a web URL to the TA & I, developed through the semester and finished by midnight, April 26, 2021.

The last few classes will be spent presenting your projects.

6

Late policy

- All projects and homework will be turned in electronically and time-stamped.
- No makeup work will be given.
- Instead, all students have 5 days of free “late time”.
This is for the entire semester, NOT per project, and counting weekends/holidays just like any other day.
 - For projects turned in late, days will be deducted from the 5 day total (or what remains of it) by the # of days late.
 - Deductions are in 1 day increments, rounding up
e.g. 10 minutes late = 1 day deducted.
 - Once the 5 days are used up, assignments will be penalized 10% / day late (rounding up), e.g., a 50 point assignment turned in 1 ½ days late would be penalized 20%, or 10 points.

7

Online homework will be via *Rosalind*: <http://rosalind.info/fag/>

Enroll specifically for BCH394P/364C at:

<http://rosalind.info/classes/enroll/8e60607640/>

Rosalind About ▾ Problems ▾ Statistics ▾ Glossary search f t My Classes ▾ edward.marcotte Log out

BCH394P/364C (Spring 2021) Systems Biology/Bioinformatics

[Edit class info](#) [Edit problems](#) [Enroll link](#) [Grade sheet](#) [Assistants](#) [Print all problems](#) [Announcements](#) [All classes](#) [Donate](#)

by Edward Marcotte at University of Texas at Austin


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.



Num	Title	Solved By	Cost	Due Date	Questions	Solutions
1	Installing Python	0	2	Jan. 27, 2021	🔒	🔒
2	Variables and Some Arithmetic	0	2	Jan. 27, 2021	🔒	🔒
3	Strings and Lists	0	2	Jan. 27, 2021	🔒	🔒
4	Conditions and Loops	0	2	Jan. 27, 2021	🔒	🔒
5	Working with Files	0	2	Jan. 27, 2021	🔒	🔒
			10			

[Found a typo?](#) [Suggest a new problem](#) [Take a tour](#)

The first homework will be due (in Rosalind) by midnight, Jan 27.

8


[About](#)
[Problems](#)
[Statistics](#)
[Glossary](#)



[My Classes](#)
[edward.marcotte](#)
[Log out](#)

Installing Python

Problem 1 @ BCH394P/364C (Spring 2021) Systems Biology/Bioinformatics ↗

Dec. 7, 2012, 12:42 p.m. by [Rosalind Team](#) Topics: [Introductory Exercises](#), [Programming](#)

Why Python? [click to expand](#)

Problem

After downloading and installing [Python](#), type `import this` into the Python command line and see what happens. Then, click the "Download dataset" button below and copy the Zen of Python into the space provided.


Time limit You'll have 5 minutes to upload the answer.



[Download dataset](#) You may make an unlimited number of attempts without being penalized.

[Questions](#)

[Found a typo?](#) [Suggest a new problem](#) [Take a tour](#)

9


[About](#)
[Problems](#)
[Statistics](#)
[Glossary](#)



[My Classes](#)
[edward.marcotte](#)
[Log out](#)

Installing Python

Problem 1 @ BCH394P/364C (Spring 2021) Systems Biology/Bioinformatics ↗

Dec. 7, 2012, 12:42 p.m. by [Rosalind Team](#) Topics: [Introductory Exercises](#), [Programming](#)

Why Python? [click to collapse](#)

Rosalind problems can be solved using any programming language. Our language of choice is [Python](#). Why? Because it's simple, powerful, and even funny. You'll see what we mean.

If you don't already have [Python](#) software, please [download and install the appropriate version for your platform](#) (Windows, Linux or Mac OS X). Please install [Python 2.x](#) (not 3.x) — it has more libraries support and many well-written guides.

After completing installation, launch [IDLE](#) (default [Python](#) development environment; it's usually installed with [Python](#), however you may need to install it separately on Linux).

You'll see a window containing three arrows, like so:

```
>>>
```

The three arrows are [Python](#)'s way of saying that it is ready to serve your every need. You are in interactive mode, meaning that any command you type will run immediately. Try typing `1+1` and see what happens.

Of course, to become a Rosalind pro, you will need to write programs having more than one line. So select **File** → **New Window** from the [IDLE](#) menu. You can now type code as you would into a text editor. For example, type the following:

```
print "Hello, World!"
```

Select **File** → **Save** to save your creation with an appropriate name (e.g., `hello.py`).

To run your program, select **Run** → **Run Module**. You'll see the result in the interactive mode window ([Python Shell](#)).

Congratulations! You just ran your first program in [Python](#)!

Problem

After downloading and installing [Python](#), type `import this` into the Python command line and see what happens. Then, click the "Download dataset" button below and copy the Zen of Python into the space provided.

Time limit You'll have 5 minutes to upload the answer.

[Download dataset](#) You may make an unlimited number of attempts without being penalized.

[Questions](#)

10

...there are quite a few good bioinformatics problems in the archives.

Rosalind is a platform for learning bioinformatics and programming through problem solving. [Take a tour](#) to get the hang of how Rosalind works.

Last win: [Hydrotelluride](#) vs. "Constructing a De Bruijn Graph", 7 minutes ago

Problems: 285 (total), users: 66794, attempts: 1014640, correct: 567494

ID	Title	Solved By	Correct Ratio	Questions	Solutions	Explanation
DNA	Counting DNA Nucleotides	35250	<div><div></div></div>			
RNA	Transcribing DNA into RNA	31498	<div><div></div></div>			
REVC	Complementing a Strand of DNA	28531	<div><div></div></div>			
FIB	Rabbits and Recurrence Relations	16249	<div><div></div></div>			
GC	Computing GC Content	16729	<div><div></div></div>			
HMM	Counting Point Mutations	18923	<div><div></div></div>			
IPRB	Mendel's First Law	10908	<div><div></div></div>			
PROT	Translating RNA into Protein	14743	<div><div></div></div>			
SUBS	Finding a Motif in DNA	15115	<div><div></div></div>			
CONS	Consensus and Profile	8423	<div><div></div></div>			
FIBD	Mortal Fibonacci Rabbits	7001	<div><div></div></div>			
GRPH	Overlap Graphs	6963	<div><div></div></div>			
IEV	Calculating Expected Offspring	6357	<div><div></div></div>			
LCSM	Finding a Shared Motif	5909	<div><div></div></div>			
LUA	Independent Alleles	3367	<div><div></div></div>			
MPRT	Finding a Protein Motif	3708	<div><div></div></div>			
MRNA	Inferring mRNA from Protein	5698	<div><div></div></div>			
ORF	Open Reading Frames	4399	<div><div></div></div>			
PERM	Enumerating Gene Orders	7860	<div><div></div></div>			
PRIM	Calculating Protein Mass	7255	<div><div></div></div>			
REVP	Locating Restriction Sites	4694	<div><div></div></div>			
SPLC	RNA Splicing	5193	<div><div></div></div>			
LEXF	Enumerating k-mers Lexicographically	4383	<div><div></div></div>			
LGIS	Longest Increasing Subsequence	1924	<div><div></div></div>			
LONG	Genome Assembly as Shortest Superstring	2195	<div><div></div></div>			
PMCH	Perfect Matchings and RNA Secondary Structures	2062	<div><div></div></div>			
PPER	Partial Permutations	2894	<div><div></div></div>			
PROB	Introduction to Random Strings	2825	<div><div></div></div>			

11

Expectations on working together

Students are welcome to discuss ideas and problems with each other, but **all programs, Rosalind homework, problem sets, and written solutions should be performed independently,**

→ except the final presentation.

tl;dr: study/discuss together
do your own programming/writing/project
collaborate on the final presentation

12



THE UNIVERSITY OF TEXAS AT AUSTIN

Student Judicial Services

Office of the Dean of Students

What is Academic Dishonesty?

In promoting a high standard of academic integrity, the University broadly defines academic dishonesty—basically, all conduct that violates this standard, including *any act designed to give an unfair or undeserved academic advantage*, such as:

- Cheating
- [Plagiarism](#)
- [Unauthorized Collaboration / Collusion](#)
- Falsifying Academic Records
- Misrepresenting Facts (e.g., providing false information to postpone an exam, obtain an extended deadline for an assignment, or even gain an unearned financial benefit)
- Any other acts (or attempted acts) that violate the basic standard of academic integrity (e.g., [multiple submissions](#)—submitting essentially the same written assignment for two courses without authorization to do so)

<https://deanofstudents.utexas.edu/conduct/academicintegrity.php>

13

- 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 (e.g. CourseHero) 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.

14

- Any materials found online (e.g. CourseHero) that are associated with you, or any suspected unauthorized sharing of materials, will be reported to Student Conduct and Academic Integrity in the Office of the Dean of Students. These reports can result in sanctions, including failure in the course.

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

15

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

16



THE UNIVERSITY OF TEXAS AT AUSTIN

Student Judicial Services

Office of the Dean of Students

Consequences of Academic Dishonesty Can Be Severe!

You may see or hear of other students engaging in some form of academic dishonesty. If so, do not assume that this misconduct is tolerated. Such violations are, in fact, regarded very seriously, often resulting in severe consequences.

Grade-related penalties are routinely assessed ("F" in the course is not uncommon), but students can also be suspended or even permanently expelled from the University for scholastic dishonesty.

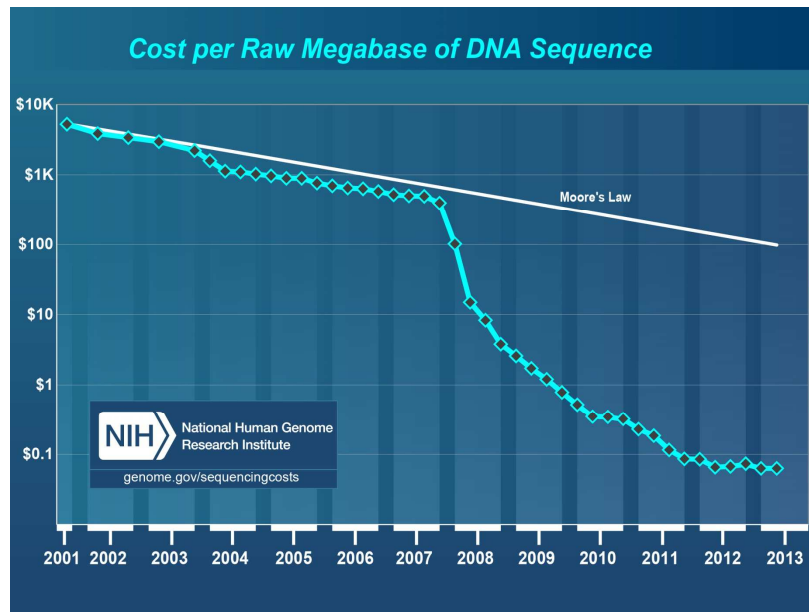
http://deanofstudents.utexas.edu/sjs/acadint_conseq.php

17

Why are we here? (practically, not existentially)

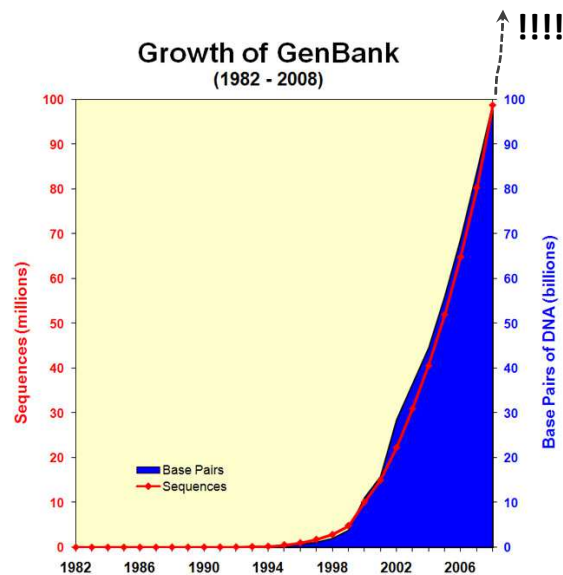
18

Pales beside the phenomenal drop in DNA sequencing costs...



21

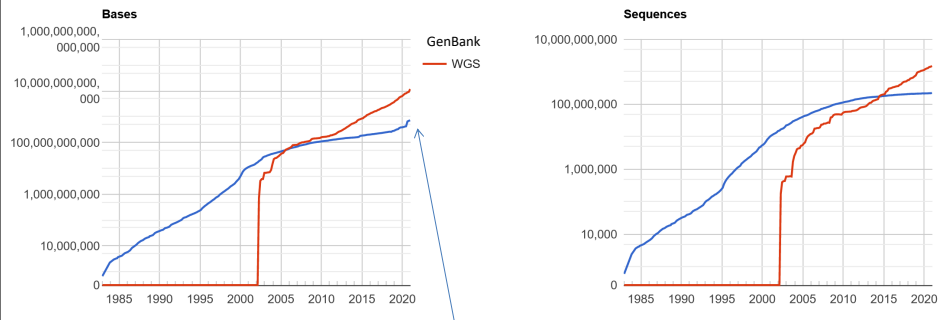
& the corresponding explosion of DNA sequencing data...



<http://www.ncbi.nlm.nih.gov/genbank/genbankstats-2008/>
<ftp://ftp.ncbi.nih.gov/genbank/gbrel.txt>

22

& the corresponding explosion of DNA sequencing data...



Here are the latest statistics...

December 2020:

723 billion bp
+
11.8 trillion bp DNA
whole genome
shotgun sequencing

Which basically
means GenBank is
falling behind
more every year!

<http://www.ncbi.nlm.nih.gov/genbank/statistics>

23

We have no choice!

Biologists are now faced with a staggering deluge of data, growing at exponential rates.

Bioinformatics offers tools and approaches to understand these data and work productively, and to build algorithmic models that help us better understand biological systems.

We'll learn some of the important basic concepts in this field, along with getting exposed to key technologies driving the field forward.

24

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

25

GENOMES, PROTEOMES, & "BIG BIOLOGY"

Gene finding algorithms

Genome assembly & how the human genome was sequenced

An introduction to large gene expression data sets

Promoter and motif finding, Gibbs sampling

Clustering algorithms, hierarchical, k-means, self-organizing maps,
force-directed maps

Classification algorithms

Principal component analysis and data transformations

NETWORK & SYNTHETIC BIOLOGY

Biological networks: metabolic, signaling, graphs, regulatory

Deep homology and the evolution of traits

Designing, simulating, and building gene circuits

Genome design and synthesis

26

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" (on zoom) demo in-class
of nanopore sequencing....**

THE FINAL COURSE PROJECT IS DUE by midnight, April 26, 2021

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