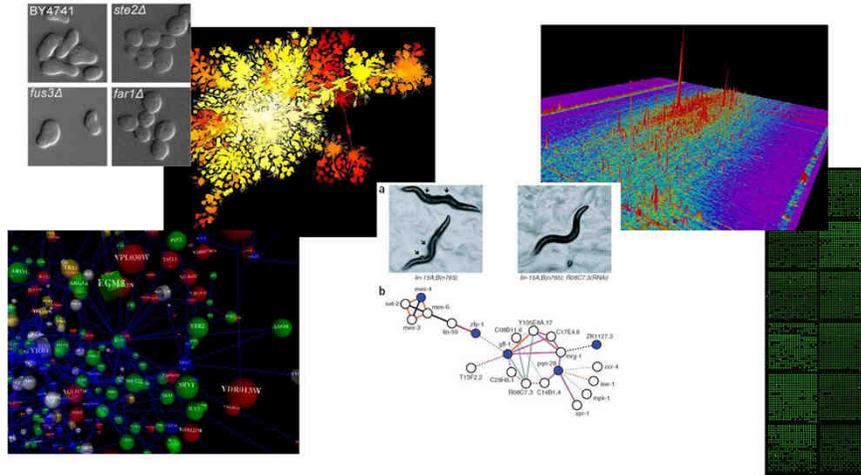


**BCH339N Systems Biology/Bioinformatics**  
(course # 54040)  
**Spring 2016 Tues/Thurs 11 – 12:30 PM BUR 212**



**Instructor: Prof. Edward Marcotte**  
**Office hours: Mon 4 PM – 5 PM**

**marcotte@icmb.utexas.edu**  
**MBB 3. 148BA**

**TA: Claire McWhite**  
**Office hours: Wed/Thurs 4 – 5 PM**  
**Phone: 512-232-3919**

**claire.mcwhite@utexas.edu**  
**MBB 3.128A**

**Probably the most important slide today!**

Course web page:

**[http://www.marcottelab.org/index.php/BCH339N\\_2016](http://www.marcottelab.org/index.php/BCH339N_2016)**

Open to biochemistry majors.

Prerequisites: Biochemistry 339F or Chemistry 339K with a grade of at least C-.

Requires basic familiarity with molecular biology & basic statistics, although varied backgrounds are expected.

**Note that this is an UNDERGRADUATE class. There is a different version intended for graduate students in alternate years (CH391L).**

**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 a few), 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 from \$26.85)

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 beginners:**

<http://www.codecademy.com/tracks/python>

## 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 3 days of class (25% of final grade)

The course project will be focused on a specific gene & will involve bioinformatics research (e.g. calculation, programming, database analysis, etc.) developed over the semester in 5 mini-assignments (4% each) 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 27, 2016.**

**The last three classes will be spent presenting your projects to each other.**

## 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.

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

Enroll specifically for BCH339N at:

<http://rosalind.info/classes/enroll/c5be9c4629>

 ROSALIND About ▾ Problems ▾ Statistics ▾ Glossary search   My Classes ▾ edward.marcotte Log out

### BCH339N Systems Biology/Bioinformatics

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

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	<a href="#">Installing Python</a>	0	2	Jan 26, 2016	⊕	⊕
2	<a href="#">Variables and Some Arithmetic</a>	0	2	Jan 26, 2016	⊕	⊕
3	<a href="#">Strings and Lists</a>	0	2	Jan 26, 2016	⊕	⊕
4	<a href="#">Conditions and Loops</a>	0	2	Jan 26, 2016	⊕	⊕
5	<a href="#">Working with Files</a>	0	2	Jan 26, 2016	⊕	⊕
			10			

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

**The first homework will be due (in Rosalind) by midnight, Jan 26.**

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

# Installing Python

## Problem 1 @ BCH339N Systems Biology/Bioinformatics

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

**Why Python?** [click to collapse](#)

**Problem**

After downloading and installing Python, type `!export` 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](#)

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

# Installing Python

## Problem 1 @ BCH339N 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 of version 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` 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 `!export` 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](#)

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

**Rosalind** About Problems Statistics Glossary search   My Classes edward marcotte Log out

**Problems** Bioinformatics Stronghold

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: [Matia Zevo](#) vs. ["Implement DistanceBetweenPatternsAndStrings"](#), 1 minute ago Problems: 254 (total), users: 33563, attempts: 568993, correct: 327689

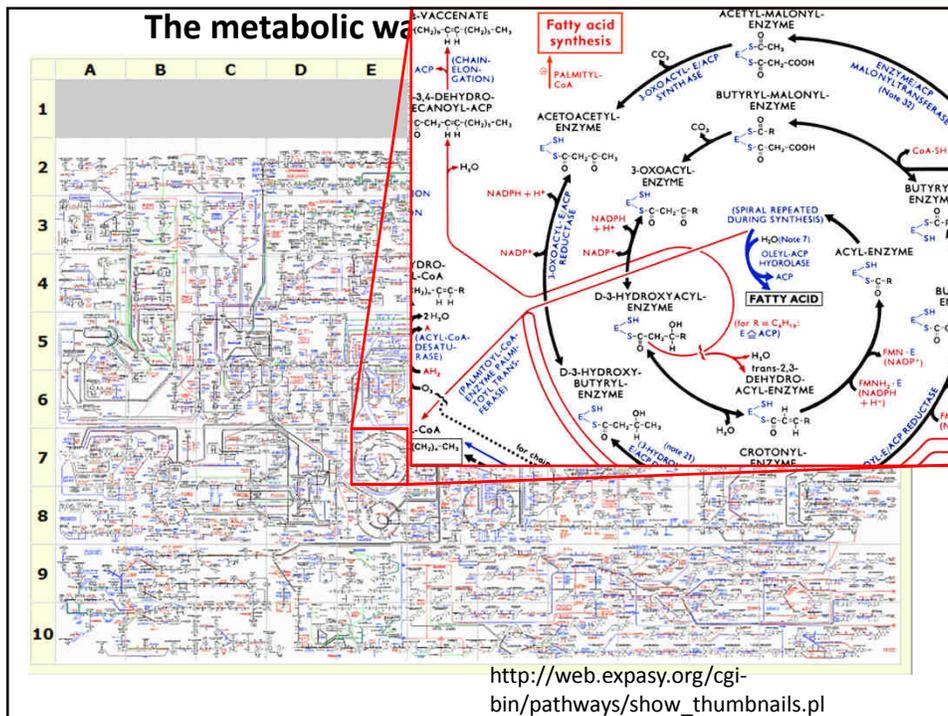
ID	Title	Solved By	Correct Ratio	Questions	Solutions	Explanation
DNA	<a href="#">Counting DNA Nucleotides</a>	20236	<div style="width: 100%; height: 10px; background-color: green;"></div>			
RNA	<a href="#">Transcribing DNA into RNA</a>	18952	<div style="width: 100%; height: 10px; background-color: green;"></div>			
REVC	<a href="#">Complementing a Strand of DNA</a>	16417	<div style="width: 100%; height: 10px; background-color: green;"></div>			
FIB	<a href="#">Rabbits and Recurrence Relations</a>	9005	<div style="width: 100%; height: 10px; background-color: green;"></div>			
GC	<a href="#">Computing GC Content</a>	9889	<div style="width: 100%; height: 10px; background-color: green;"></div>			
HAMM	<a href="#">Counting Point Mutations</a>	11231	<div style="width: 100%; height: 10px; background-color: green;"></div>			
IPRB	<a href="#">Mendel's First Law</a>	5970	<div style="width: 100%; height: 10px; background-color: green;"></div>			
PROT	<a href="#">Translating RNA into Protein</a>	8511	<div style="width: 100%; height: 10px; background-color: green;"></div>			
SUBS	<a href="#">Finding a Motif in DNA</a>	8958	<div style="width: 100%; height: 10px; background-color: green;"></div>			
CONS	<a href="#">Consensus and Profile</a>	5997	<div style="width: 100%; height: 10px; background-color: green;"></div>			
FIBD	<a href="#">Morris Fibonacci Rabbits</a>	4034	<div style="width: 100%; height: 10px; background-color: green;"></div>			
GRPH	<a href="#">Overlap Graphs</a>	4264	<div style="width: 100%; height: 10px; background-color: green;"></div>			
IEV	<a href="#">Calculating Expected Offspring</a>	3691	<div style="width: 100%; height: 10px; background-color: green;"></div>			
LCSM	<a href="#">Finding a Shared Motif</a>	3586	<div style="width: 100%; height: 10px; background-color: green;"></div>			
LIA	<a href="#">Independent Alleles</a>	1917	<div style="width: 100%; height: 10px; background-color: green;"></div>			
MPRT	<a href="#">Finding a Protein Motif</a>	2172	<div style="width: 100%; height: 10px; background-color: green;"></div>			
MRNA	<a href="#">Inferring mRNA from Protein</a>	3435	<div style="width: 100%; height: 10px; background-color: green;"></div>			
ORF	<a href="#">Open Reading Frames</a>	2620	<div style="width: 100%; height: 10px; background-color: green;"></div>			
PERM	<a href="#">Enumerating Gene Orders</a>	5001	<div style="width: 100%; height: 10px; background-color: green;"></div>			
PRTM	<a href="#">Calculating Protein Mass</a>	4278	<div style="width: 100%; height: 10px; background-color: green;"></div>			
REVP	<a href="#">Locating Restriction Sites</a>	2945	<div style="width: 100%; height: 10px; background-color: green;"></div>			
SPLC	<a href="#">RNA Splicing</a>	3691	<div style="width: 100%; height: 10px; background-color: green;"></div>			
LEXF	<a href="#">Enumerating k-mers Lexicographically</a>	2884	<div style="width: 100%; height: 10px; background-color: green;"></div>			
LGIS	<a href="#">Longest Increasing Subsequence</a>	1159	<div style="width: 100%; height: 10px; background-color: green;"></div>			
LONG	<a href="#">Genome Assembly as Shortest Superstring</a>	1358	<div style="width: 100%; height: 10px; background-color: green;"></div>			
PMCH	<a href="#">Perfect Matchings and RNA Secondary Structures</a>	1214	<div style="width: 100%; height: 10px; background-color: green;"></div>			
PPER	<a href="#">Partial Permutations</a>	1803	<div style="width: 100%; height: 10px; background-color: green;"></div>			
PROB	<a href="#">Introduction to Random Strings</a>	1729	<div style="width: 100%; height: 10px; background-color: green;"></div>			

## Expectations on working together

Students are welcome to discuss ideas and problems with each other, but **all programs, Rosalind homework, 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

Why are we here? (practically, not existentially)

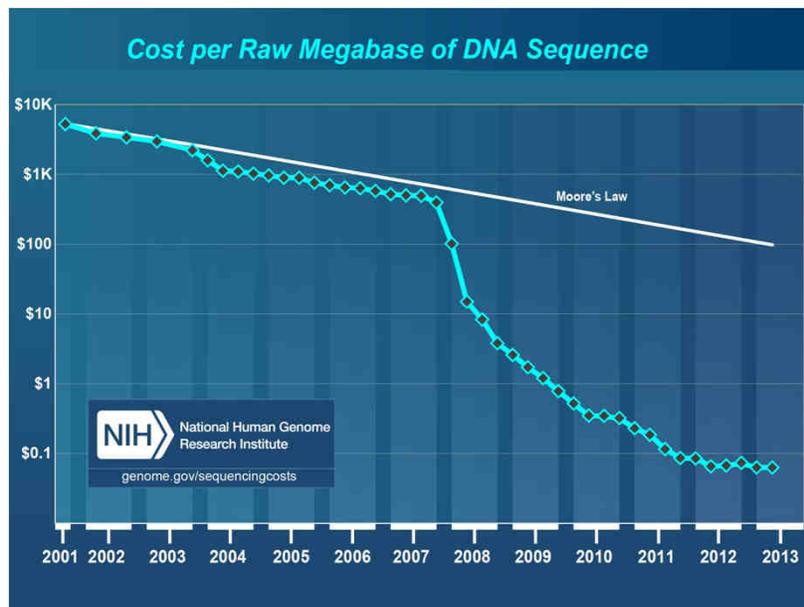


## Our current knowledge of human metabolism...

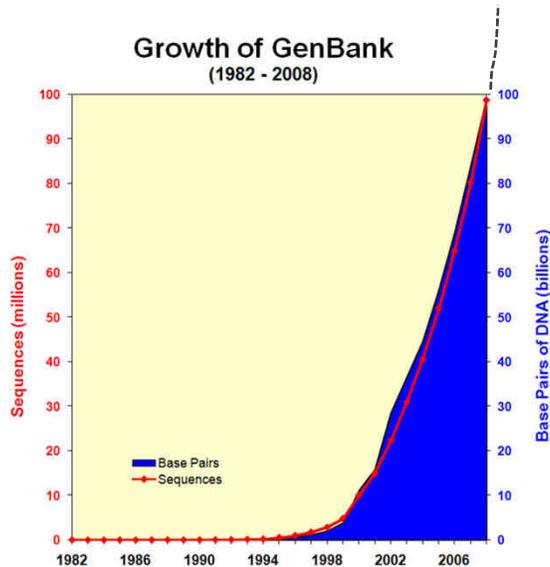
Total number of reactions	7,440
Total number of metabolites	5,063
Number of unique metabolites	2,626
Number of metabolites in extracellular space	642
Number of metabolites in cytoplasm	1,878
Number of metabolites in mitochondrion	754
Number of metabolites in nucleus	165
Number of metabolites in endoplasmic reticulum	570
Number of metabolites in peroxisome	435
Number of metabolites in lysosome	302
Number of metabolites in Golgi apparatus	317
Number of transcripts	2,194
Number of unique genes	1,789

Nat Biotechnol. 2013 May;31(5):419-25

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

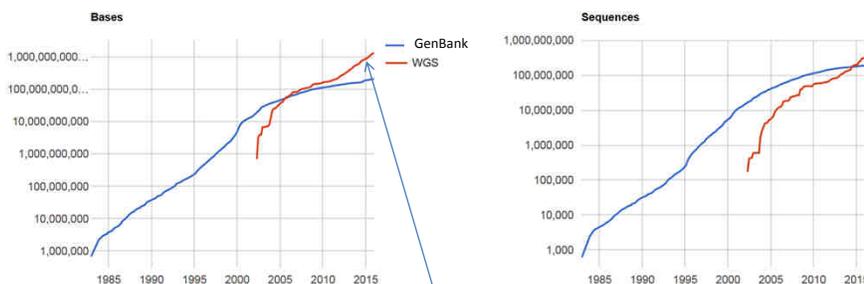


**& the corresponding explosion of DNA sequencing data...**



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

**& the corresponding explosion of DNA sequencing data...**



Here are the latest statistics...

**December 2015:**  
 203 billion bp  
 +  
 1.3 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>

**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.**

### **Specifically...**

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

#### **BASICS OF PROGRAMMING**

Introduction to Rosalind

A Python programming primer for non-programmers

#### **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

### **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  
Classifiers, k-nearest neighbors, Mahalanobis distance  
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

Homology, orthology, and evolutionary trees  
3D modeling of protein structures  
Next- (& next-next-) generation DNA and RNA sequencing  
Mass spectrometry shotgun proteomics  
Genome engineering

**\*\*\* THE FINAL GENE PAGES PROJECT IS DUE by midnight, April 27, 2016 \*\*\***

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