

# CSE 527

## Computational Biology

<http://www.cs.washington.edu/527>

### Lecture 1: Overview & Bio Review

Autumn 2004

Larry Ruzzo

#### CSE 527: Computational Biology, Autumn 2004

An introduction to the use of computational methods for the understanding of biological systems at the molecular level. Intended for graduate students in biological sciences interested in learning about algorithms and computational methods, and for graduate students in computer science, mathematics or statistics interested in applications of these fields to molecular biology.

**Time:** MW 12:00-1:20

**Place:** EE1 025

**Instructor:** Larry Ruzzo (CSE 554, ruzzo@cs.washington.edu)

**Course web page:** <http://www.cs.washington.edu/527>

**Course mailing list:** [cse527@cs.washington.edu](mailto:cse527@cs.washington.edu)

**Catalog Description (somewhat out of date):** CSE 527 Computational Biology (3) introduces computational methods for understanding biological systems at the molecular level. Problem areas such as mapping and sequencing, sequence analysis, structure prediction, phylogenetic inference, regulatory analysis. Techniques such as dynamic programming, Markov models, expectation-maximization, local search. Prerequisite: graduate standing in biological, computer, mathematical or statistical science, or permission of instructor.

**Workload:** Notes, problem sets, project. We encourage projects in which a biologist and a mathematical scientist collaborate to model/solve a biological problem.

**Desired Prerequisites:** Ideally, students will have a considerable knowledge of one of computer science, biology, or probability/statistics, plus introductory knowledge of the other two. We'll try to supplement as needed (via lecture, outside reading, project teams, etc.) so that everyone has enough background in the immediately relevant areas to fruitfully proceed.

#### Rough Course Outline

I am unlikely to have time to cover all of this. If you have particular interests, let me know and I'll try to prioritize the in-demand topics.

#### Essential Background from Molecular Biology

**Sequence Analysis:** Statistical modeling of families of DNA or protein sequences; profiles, motif discovery, hidden Markov Models, Expectation-Maximization algorithm, Gibbs sampling, Gene finding.

**Molecular Structure Prediction (time permitting):** RNA secondary structure prediction, SCFGs and covariance models; the protein folding problem; protein threading.

**Microarray Analysis:** Clustering, classification, feature selection for analysis of large scale gene expression data sets generated by microarrays and similar technologies.

He who asks is a fool for five minutes, but he who does not ask remains a fool forever.

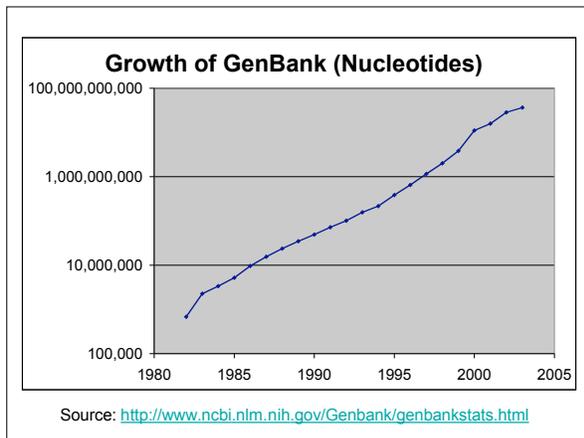
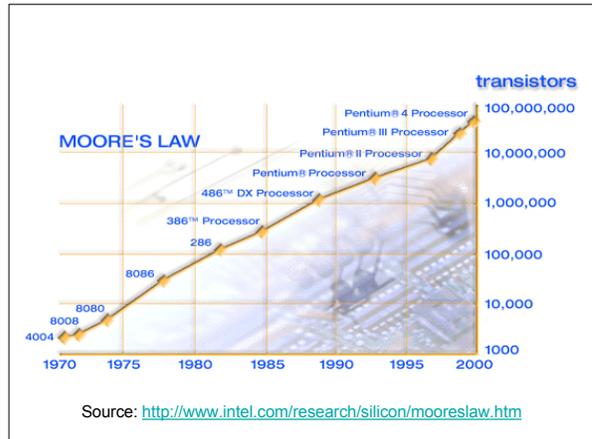
-- Chinese Proverb

## Related Courses

- Genome 540/541 (Winter/Spring)
  - Intro. To Comp. Mol. Bio.
- Stat/Biostat 578 (A 2004)
  - Statistical Analysis of Microarrays
- CSE590CB (AWS)
  - Reading & Research in Comp. Bio.
  - Monday's, 3:30 (MEB 243 this quarter)
  - <http://www.cs.washington.edu/590cb>
- Combi Seminar (Genome 521; AWS)
  - Wednesday's 1:30 K069 (sometimes 3:30 Hitch 132)

## Homework #1

- Find & read a good primer on “bio for cs” (or vice versa, as appropriate) e.g., see ones listed on 590cb page
- Email me a few sentences saying
  - What you read (give me a link or citation)
  - Critique it for your meeting your needs
  - Who would it have been good for, if not you



## What's all the fuss?

- The human genome is “finished”...
- Even if it were, that's only the beginning
- Explosive growth in biological data is revolutionizing biology & medicine

“All pre-genomic lab techniques are obsolete”

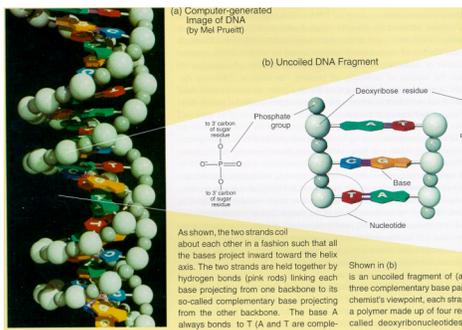
(and computation and mathematics are crucial to post-genomic analysis)

## A *VERY* Quick Intro To Molecular Biology

## The Genome

- The hereditary info present in every cell
- DNA molecule -- a long sequence of *nucleotides* (A, C, T, G)
- Human genome -- about  $3 \times 10^9$  nucleotides
- The genome project -- extract & interpret genomic information, apply to genetics of disease, better understand evolution, ...

## The Double Helix



## DNA

- Discovered 1869
- Role as carrier of genetic information - much later
- The Double Helix - Watson & Crick 1953
- Complementarity
  - $A \leftrightarrow T$      $C \leftrightarrow G$

## Genetics - the study of heredity

- A *gene* -- classically, an abstract heritable attribute existing in variant forms (*alleles*)
- *Genotype vs phenotype*
- Mendel
  - Each individual two copies of each gene
  - Each parent contributes one (randomly)
  - Independent assortment

## Cells

- Chemicals inside a sac - a fatty layer called the *plasma membrane*
- *Prokaryotes* (e.g., bacteria) - little recognizable substructure
- *Eukaryotes* (all multicellular organisms, and many single celled ones, like yeast) - genetic material in nucleus, other organelles for other specialized functions

## Chromosomes

- 1 pair of DNA molecules (+ protein wrapper)
- Most prokaryotes have just 1 chromosome
- Eukaryotes - all cells have same number of chromosomes, e.g. fruit flies 8, humans & bats 46, rhinoceros 84, ...

## Mitosis/Meiosis

- Most “higher” eukaryotes are *diploid* - have homologous pairs of chromosomes, one maternal, other paternal (exception: sex chromosomes)
- *Mitosis* - cell division, duplicate each chromosome, 1 copy to each daughter cell
- *Meiosis* - 2 divisions form 4 *haploid* gametes (egg/sperm)
  - *Recombination/crossover* -- exchange maternal/paternal segments

## Proteins

- Chain of amino acids, of 20 kinds
- Proteins are the major functional elements in cells
  - Structural
  - Enzymes (catalyze chemical reactions)
  - Receptors (for hormones, other signaling molecules, odors,....)
  - Transcription factors
  - ...
- 3-D Structure is crucial: the protein folding problem

## The “Central Dogma”

- Genes encode proteins
- DNA transcribed into messenger RNA
- RNA translated into proteins
- Triplet code (codons)

## The Genetic Code

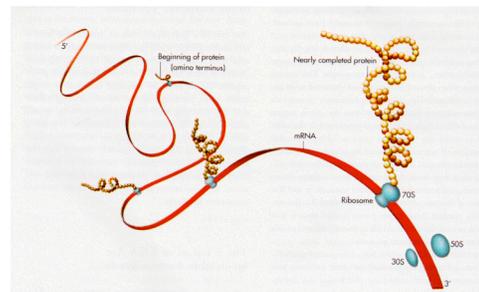
(a) RNA Codons for the Twenty Amino Acids

		Second base				
		U	C	A	G	
First base	U	Phe Phe Leu Leu	Ser Ser Ser Ser	Tyr Tyr STOP STOP	Cys Cys STOP TTP	U C A G
	C	Leu Leu Leu Leu	Pro Pro Pro Pro	His His Gln Gln	Arg Arg Arg Arg	U C A G
	A	Ile Ile Ile Met (start)	Thr Thr Thr Thr	Asn Asn Lys Lys	Ser Ser Arg Arg	U C A G
G	Val Val Val Val	Ala Ala Ala Ala	Asp Asp Glu Glu	Gly Gly Gly Gly	U C A G	

**Amino-acid abbreviations**

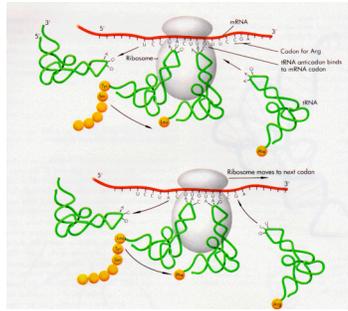
Ala = Alanine  
 Arg = Arginine  
 Asp = Aspartic acid  
 Asn = Asparagine  
 Cys = Cysteine  
 Glu = Glutamic acid  
 Gln = Glutamine  
 Gly = Glycine  
 His = Histidine  
 Ile = Isoleucine  
 Leu = Leucine  
 Lys = Lysine  
 Met = Methionine  
 Phe = Phenylalanine  
 Pro = Proline  
 Ser = Serine  
 Thr = Threonine  
 Trp = Tryptophan  
 Tyr = Tyrosine  
 Val = Valine

## Translation: mRNA → Protein



Watson, Gilman, Wilkowitz, & Zoller, 1992

## Ribosomes



Watson, Gilman, Wilkowsk, & Zoller, 1992

## Gene Structure

- Transcribed 5' to 3'
- Promoter region and transcription factor binding sites precede 5'
- Transcribed region includes 5' and 3' untranslated regions
- In eukaryotes, most genes also include introns, spliced out before export from nucleus, hence before translation

## Genome Sizes

	Base Pairs	Genes
<i>Mycoplasma genitalium</i>	580,073	483
<i>E. coli</i>	4,639,221	4,290
<i>Saccharomyces cerevisiae</i>	12,495,682	5,726
<i>Caenorhabditis elegans</i>	95.5 x 10 <sup>6</sup>	19,820
<i>Arabidopsis thaliana</i>	115,409,949	25,498
<i>Drosophila melanogaster</i>	122,653,977	13,472
Humans	3.3 x 10 <sup>9</sup>	~25,000

## Genome Surprises

- Humans have < 1/3 as many genes as expected
- But perhaps more proteins than expected, due to *alternative splicing*
- There are unexpectedly many *non-coding RNAs*
- Many other non-coding regions are highly conserved, e.g., across all mammals

... and much more ...

- Read one of the many intro surveys or books for much more info.