

# CSEP590A Computational Biology

<http://www.cs.washington.edu/csep590a>

Larry Ruzzo  
Autumn 2008



UW CSE Computational Biology Group

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

-- Chinese Proverb

## Tonight

Admin

Why Comp Bio?

The world's shortest Intro. to Mol. Bio.

Admin Stuff

## Course Mechanics & Grading

Web <http://www.cs.washington.edu/csep590a>

Reading

In class discussion

Homeworks

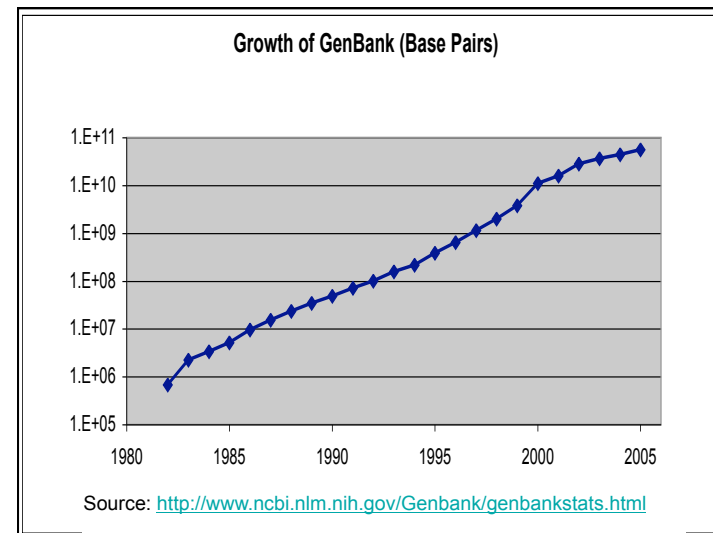
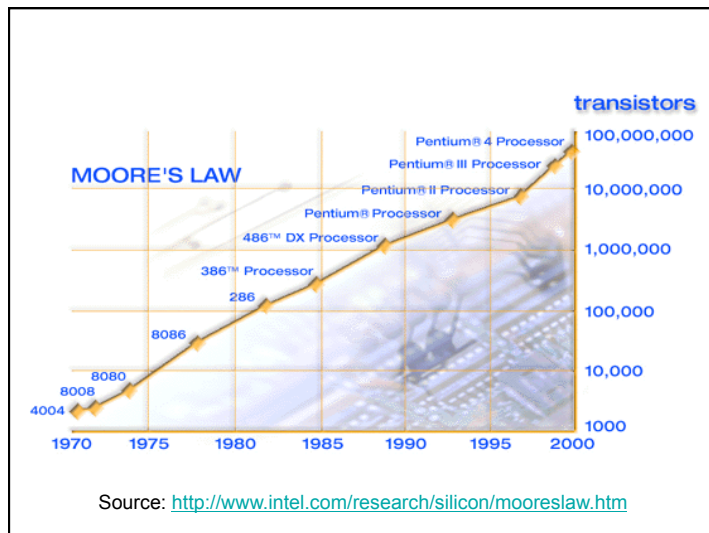
reading blogs

paper exercises

programming

No exams, but possible oversized last homework in lieu of final

## Background & Motivation



# The Human Genome Project

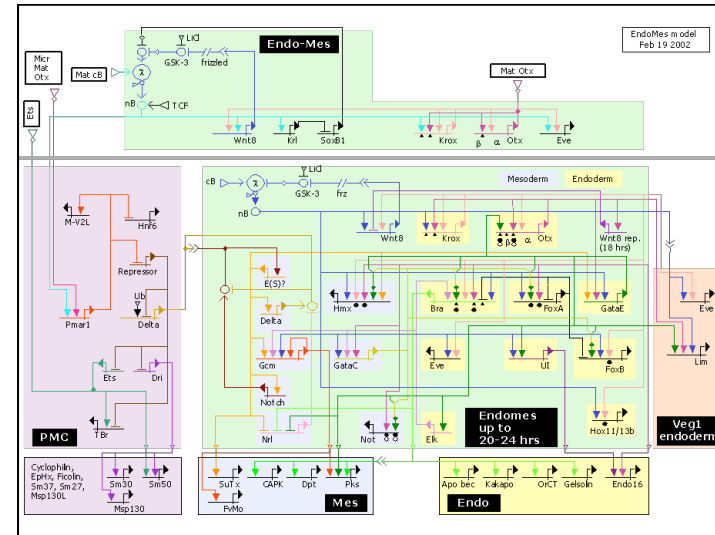
```

1 gagccccggc cgggggacgg gcggcgggat agcgggaccc cggcgcggcg gtgcgcttca
61 gggcgcagcg cgggcccagc accgagcccc gggcgcggca agaggcggcg ggagccggtg
121 gggcctcgcg atcatgcgtc gagggcgtct gctggagatc gccctgggat ttaccgtgct
181 tttagcgtcc tacacgagcc atggggcgga cgccaatttg gaggtcggga acgtgaagga
241 aaccagagcc agtcgggcca agagaagagg cggtgaggga cacgacggcg ttaaaggacc
301 caatgtctgt ggatcacggtt ataatcetta ctgtgcccct ggatggaaaa ccttacctgg
361 cggaaatcag tgtattgtcc caatttgccc gcattctctg ggggatggat tttgttcgag
421 gccaaatgat tgcacttgcc catctgggta gatagctcct tctgtggct ccagatccat
481 acaacactgc aatattcgct gtatgaatgg aggtagctgc agtgacgatc actgtctatg
541 ccagaaagga tacatagggc ctoactgtgg acaacctgtt tgtgaaagtg gctgtctcaa
601 tggaggaagg tgtgtggccc caaatcgatg tgcattgcat tacggattta ctggaccoca
661 gtgtgaaaga gattacagga caggcccatg ttttactgtg atcagcaacc agatgtgcca
721 gggacaactc agcgggatgg tctgcacaaa acagctctgc tgtgccacag tcggccgagc
781 ctggggccac cctgtgaga tgtgtcctgc ccagcctcac cctgcgccgc gtgggtccat
841 tccaaatatc cgcacgggag ctgttcaaga tgtggatgaa tgccaggcca tcccccggct
901 ctgtcaggga ggaattgca ttaatactgt tgggtctttt gactgcaaat gcctctgctg
961 acacaaactt aatgaagtgt cacaaaaatg tgaagatatt gatgaatgca gcaccattcc
1021 ...

```

## Goals

- Basic biology
- Disease diagnosis/prognosis/treatment
- Drug discovery, validation & development
- Individualized medicine
- ...



**“High-Throughput BioTech”**

- Sensors:**
  - DNA sequencing
  - Microarrays/Gene expression
  - Mass Spectrometry/Proteomics
  - Protein/protein & DNA/protein interaction
- Controls:**
  - Cloning
  - Gene knock out/knock in
  - RNAi

**Floods of data**

**“Grand Challenge” problems**

## 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)

## CS Points of Contact & Opportunities

### Scientific visualization

Gene expression patterns

### Databases

Integration of disparate, overlapping data sources

Distributed genome annotation in face of shifting underlying genomic coordinates, individual variation, ...

### AI/NLP/Text Mining

Information extraction from text with inconsistent nomenclature, indirect interactions, incomplete/inaccurate models, ...

### Machine learning

System level synthesis of cell behavior from low-level heterogeneous data (DNA seq, gene expression, protein interaction, mass spec,...)

...

### Algorithms

## An Algorithm Example: ncRNAs

The “Central Dogma”:  
DNA -> messenger RNA -> Protein

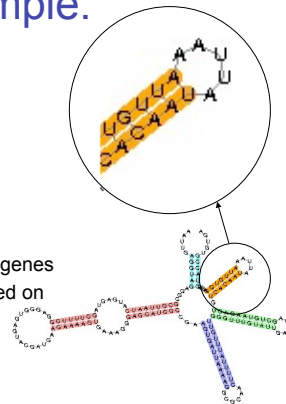
Last ~5 years: many examples  
of functionally important ncRNAs

175 -> 350 families just in last 6 mo.

Much harder to find than protein-coding genes

Main method - Covariance Models (based on  
stochastic context free grammars)

Main problem - Sloooow ...  $O(nm^4)$



## “Rigorous Filtering” - Z. Weinberg

Convert CM to HMM

(AKA: stochastic CFG to stochastic regular grammar)

Do it so HMM score *always*  $\geq$  CM score

Optimize for most aggressive filtering subject to constraint that  
score bound maintained

A large convex optimization problem

Filter genome sequence with fast HMM, run (slow) CM only on  
sequences above desired HMM threshold; guaranteed not to miss  
anything

Newer, more elaborate techniques pulling in by secondary  
structure features for further searching  
(uses automata theory, dynamic programming, Dijkstra, more  
optimization stuff,...)

details  
CENSORED  
(but stay tuned...)  
Plenty of CS here

## Results

Typically 200-fold speedup or more  
Finding dozens to hundreds of new  
ncRNA genes in many families  
Has enabled discovery of many new  
families

Newer, more elaborate techniques pulling in key secondary  
structure features for better searching  
(uses automata theory, dynamic programming, Dijkstra, more  
optimization stuff,...)

## The Mission

“Solving **Today’s** challenging  
Computer Science problems  
for **Tomorrow’s** biologists”

## More Admin

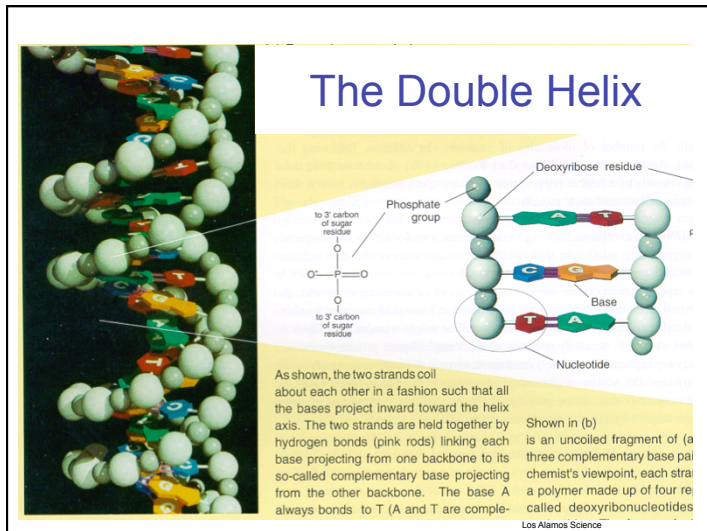
## Course Focus & Goals

Mainly sequence analysis  
Algorithms for alignment, search, & discovery  
    Specific sequences, general types (“genes”, etc.)  
    Single sequence and comparative analysis  
Techniques: HMMs, EM, MLE, Gibbs, Viterbi...  
Enough bio to motivate these problems, including very  
    light intro to modern biotech supporting them  
Math/stats/cs underpinnings thereof  
Applied to real data

## 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, ...



## DNA

Discovered 1869  
Role as carrier of genetic information -  
much later  
4 "bases":  
adenine (A), cytosine (C), guanine (G), thymine (T)  
The Double Helix - Watson & Crick 1953  
Complementarity  
 $A \leftrightarrow T \quad C \leftrightarrow G$

## 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 (complementary) 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

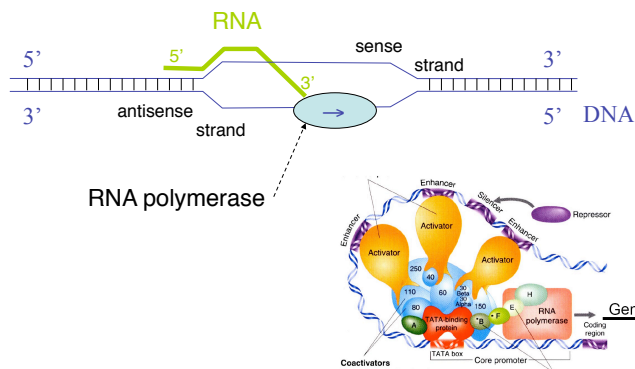
## Proteins

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

## The “Central Dogma”

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

## Transcription: DNA → RNA



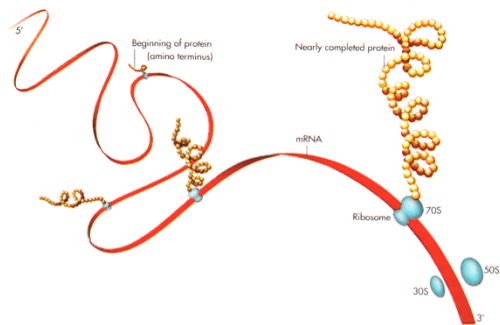
## Codons & The Genetic Code

		Second Base				
		U	C	A	G	
First Base	U	Phe Leu	Ser Ser	Tyr Stop	Cys Stop	U C A G
	C	Leu Leu Leu	Pro Pro Pro	His His Gln	Arg Arg Arg	U C A G
	A	Ile Ile Met/Start	Thr Thr Thr	Asn Asn Lys	Ser Ser 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
						Third Base

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

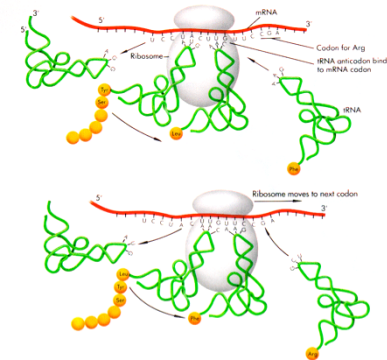


## Translation: mRNA → Protein



Watson, Gilman, Wilkowsky, & Zoller, 1992

## Ribosomes



Watson, Gilman, Wilkowsky, & Zoller, 1992

## Gene Structure

Transcribed 5' to 3'

Promoter region and transcription factor binding sites (usually) 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
MimiVirus	1,200,000	1,260
<i>E. coli</i>	4,639,221	4,290
<i>Saccharomyces cerevisiae</i>	12,495,682	5,726
<i>Caenorhabditis elegans</i>	95,500,000	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*  
-- more than protein-coding genes, by some  
estimates  
Many other non-coding regions are highly  
conserved, e.g., across all vertebrates

... and much more ...

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

## Homework #1 (partial)

Read Hunter's "bio for cs" primer;  
Find & read another  
Post 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  
See class web for more details