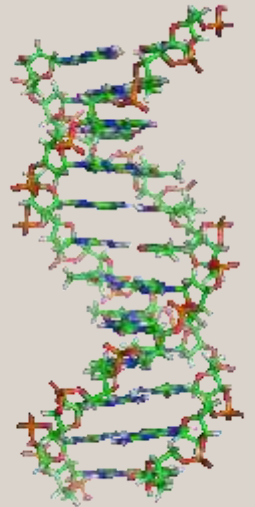# Introduction to Bioinformatics

*Dan Lopresti*

Computer Science and Engineering

Office Building C 337

dal9@lehigh.edu
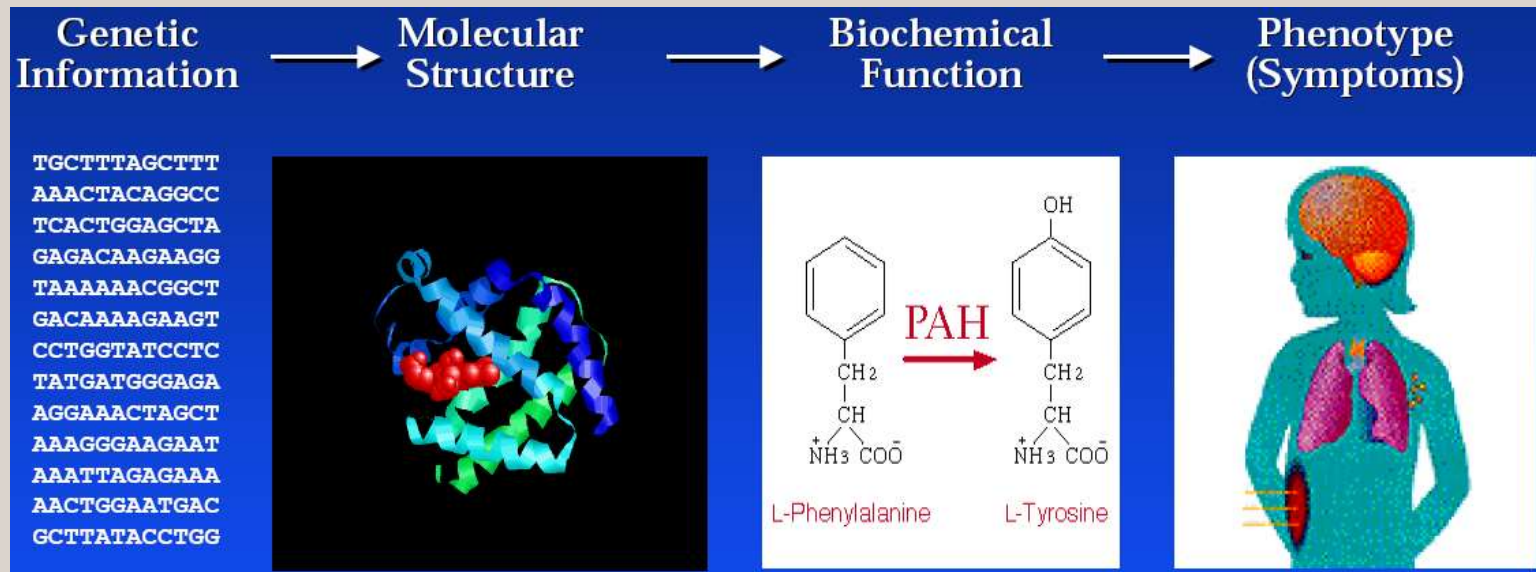
HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# In 2017 when I gave this talk ...

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Motivation

"Biology easily has 500 years of exciting problems to work on."  *Donald Knuth (famous computer scientist)*



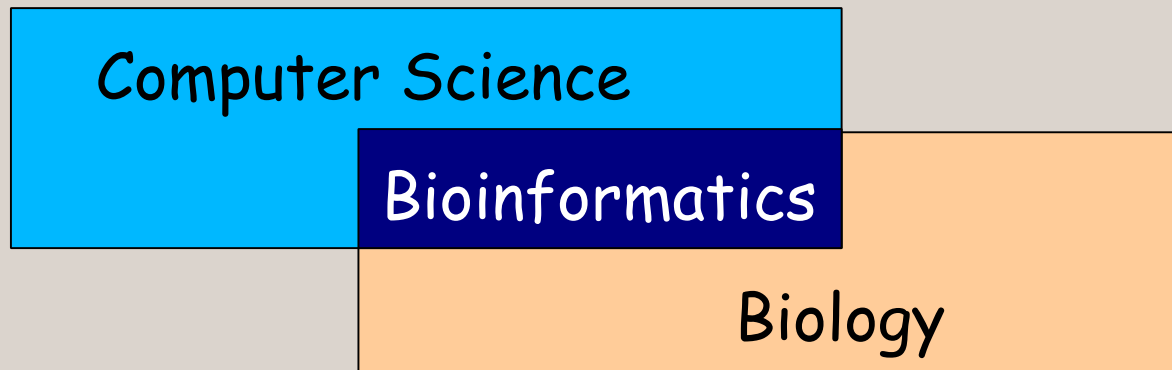By developing techniques for analyzing sequence data and structures, we can attempt to understand basis of life.

http://cmgm.stanford.edu/biochem218/

# Bioinformatics

What is bioinformatics?   Application of methods from computer science to biology.

Computer Science

Bioinformatics

Biology

Why is it interesting?
- Important problems.
- Massive quantities of data.
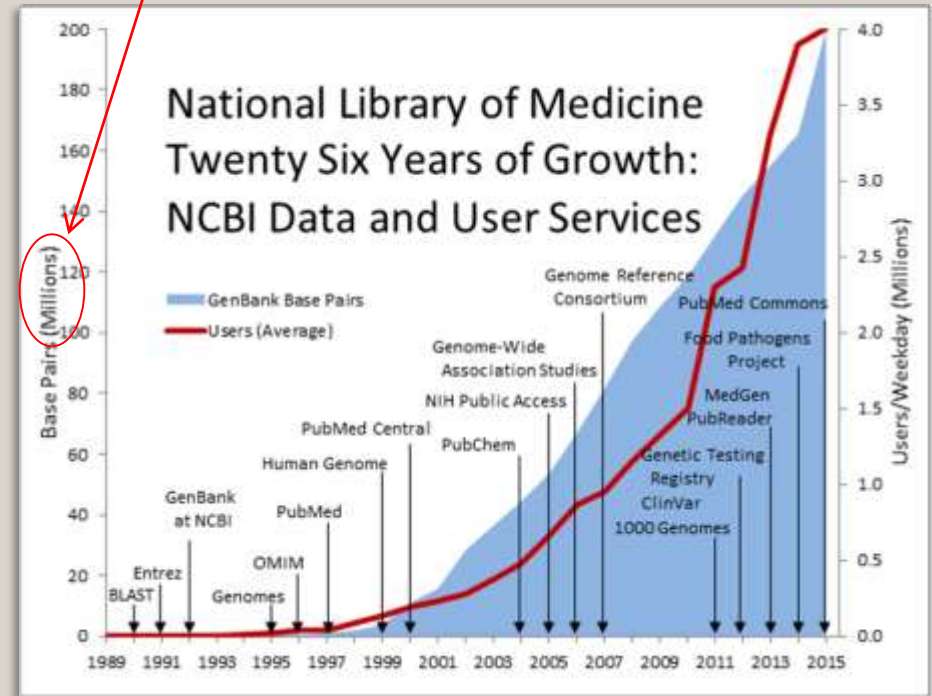- Great need for efficient solutions.
- Success is rewarded.

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Data Explosion

Our genetic identity is encoded in long molecules made up of four basic units, the nucleic acids:

    (1) *Adenine,*
    (2) *Cytosine,*
    (3) *Guanine,*
    (4) *Thymine.*

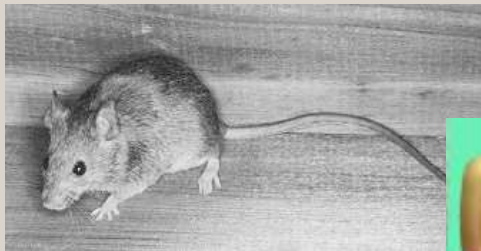To first approximation, DNA is a language over a four character alphabet, {*A, C, G, T* }.

NLM / NIH seems to have made a mistake: this should be billions, not millions!



https://www.nlm.nih.gov/about/2017CJ.html

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Genomes

Set of chromosomes that determines an organism is known as its *genome*.

Us →

Mus musculus

Poaceae

Zea mays

**GenBank Release 121.0 — December 15, 2000**

| Species | Haploid genome size | Bases | Entries |
|---|---|---|---|
| Homo sapiens | 3,400,000,000 | 6,702,881,570 | 3,918,724 |
| Mus musculus | 3,454,200,000 | 1,291,602,139 | 2,456,194 |
| Drosophila melanogaster | 180,000,000 | 487,561,384 | 166,554 |
| Arabidopsis thaliana | 100,000,000 | 242,674,129 | 181,388 |
| Caenorhabditis elegans | 100,000,000 | 203,544,197 | 114,553 |
| Tetraodon nigroviridis | 350,000,000 | 165,539,271 | 188,993 |
| Oryza sativa | | | 411 |
| Rattus norvegicus | | | 598 |
| Bos taurus | | | 473 |
| Glycine max | | | 802 |
| Medicago truncatula | | | 535 |
| Trypanosoma brucei | | | 534 |
| Lycopersicon esculent | | | 112 |
| Giardia intestinalis | | | 828 |
| Strongylocentrotus pu | | | 532 |
| Entamoeba histolytica | | | 038 |
| Hordeum vulgare | — | 44,489,692 | 57,779 |
| Danio rerio | 1,900,000,000 | 40,906,902 | 83,726 |
| Zea mays | 5,000,000,000 | 36,885,212 | 77,506 |
| Saccharomyces cerevisiae | 12,067,280 | 32,779,082 | 18,361 |

Conclusion:  size does <u>not</u> matter!
(But you already knew this.  😊 )

http://www.cbs.dtu.dk/databases/DOGS/
http://www.nsrl.ttu.edu/tmot1/mus_musc.htm
http://www.oardc.ohio-state.edu/seedid/single.asp?strID=324

# Comparative Genomics



**Mouse and Human Genetic Similarities**

How did we decipher these relationships?

Courtesy Lisa Stubbs
Oak Ridge National Laboratory

YGA 98-075R2

http://www.ornl.gov/sci/techresources/Human_Genome/graphics/slides/ttmousehuman.shtml

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
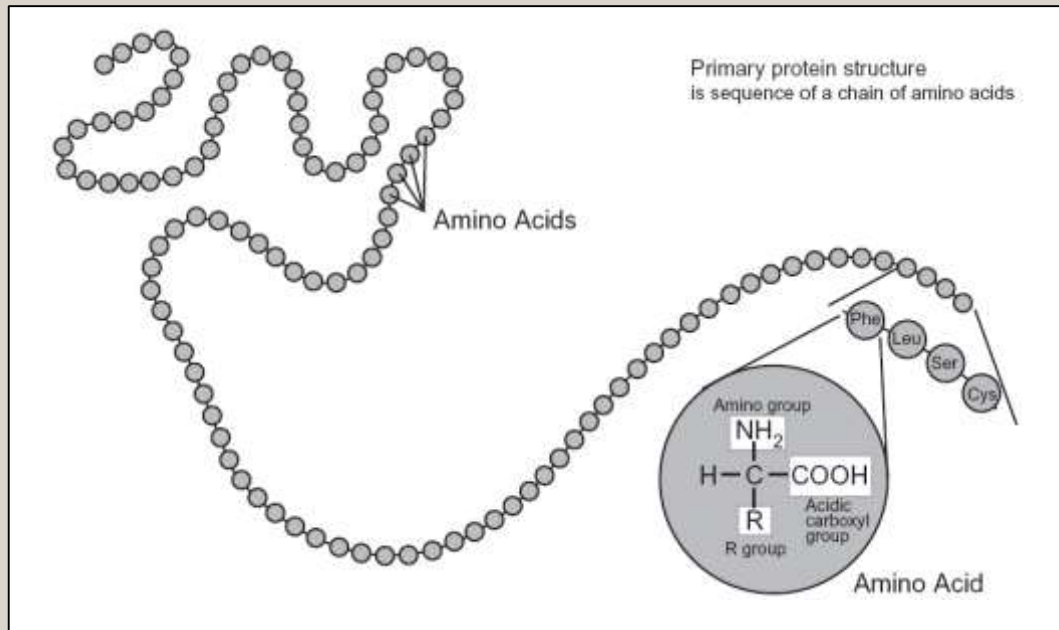CSE

# Algorithms are Central

An *algorithm* is a precisely-specified series of steps to solve a particular problem of interest.

- Develop model(s) for task at hand.

- Study inherent computational complexity:
  - Can task be phrased as an optimization problem?
  - Can it be solved efficiently?  Speed, memory, etc.
  - If we can't find good algorithm, can we prove task hard?
  - If known to be hard, is there approximation algorithm (works some of the time or comes close to optimal)?

- Conduct experimental evaluations (iterate above steps).

HHMI
Howard Hughes Medical Institute

LEHIGH
U N I V E R S I T Y.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Sequence Nature of Biology

*Macromolecules* are chains of simpler molecules.



Primary protein structure is sequence of a chain of amino acids

Amino Acids

Phe
Leu
Ser
Cys

Amino group
NH₂
H—C—COOH
R
R group
Acidic carboxyl group

Amino Acid

http://www.accessexcellence.org/AB/GG/aminoAcid.html



AUCG's
ATCG's

Nitrogenous Bases

Base pair

Sugar phosphate backbone

RNA
Ribonucleic acid

DNA
Deoxyribonucleic acid

http://www.accessexcellence.org/AB/GG/rna.html

In case of proteins, these basic building blocks are *amino acids*.

In DNA and RNA, they are *nucleotides*.

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering

CSE

# NCBI GenBank

National Center for Biotechnology Information (NCBI), a branch of National Institutes of Health (NIH), maintains *GenBank*, a worldwide repository of genetic sequence data (all publicly available DNA sequences).



Massive quantities of sequence data ⇒ need for good computational techniques.

http://www.ncbi.nlm.nih.gov/

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Reading DNA



**This is known as sequencing.**

http://www.apelex.fr/anglais/applications/sommaire2/sanger.htm
http://www.iupui.edu/~wellsctr/MMIA/htm/animations.htm

*Gel electrophoresis separates mixture of molecules in a gel media by application of an electric field.*

In general, molecules with similar lengths will migrate same distance.

Make DNA fragments that end at each base. Then run gel and read off sequence:  ATCGTG …

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Reading DNA

Original sequence: ATCGTGTCGATAGCGCT

**G**
```
ATCG
ATCGTG
ATCGTGTCG
ATCGTGTCGATAG
ATCGTGTCGATAGCG
```

**A**
```
A
ATCGTGTCGA
ATCGTGTCGATA
```

**T**
```
AT
ATCGT
ATCGTGT
ATCGTGTCGAT
ATCGTGTCGATAGCGCT
```

**C**
```
ATC
ATCGTGTC
ATCGTGTCGATAGC
ATCGTGTCGATAGCGC
```

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Sequencing a Genome

Most genomes are enormous (e.g., $10^{10}$ base pairs for human). But current sequencing technology only allows biologists to determine ~$10^3$ base pairs at a time.

Leads to some very interesting problems in bioinformatics!

*Genetic linkage map*
*($10^7$ – $10^8$ base pairs)*

*Physical map*
*($10^5$ – $10^6$ base pairs)*

*ACTAGCTGATCGATTTAGCAGCAG...*

*Sequencing*
*($10^3$ – $10^4$ base pairs)*

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Sequencing a Genome

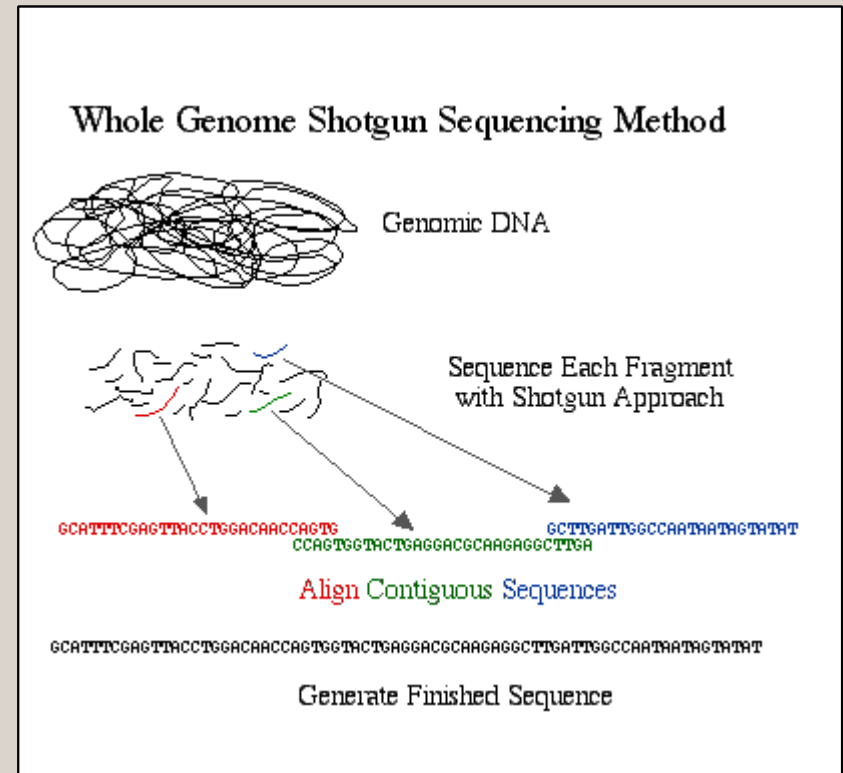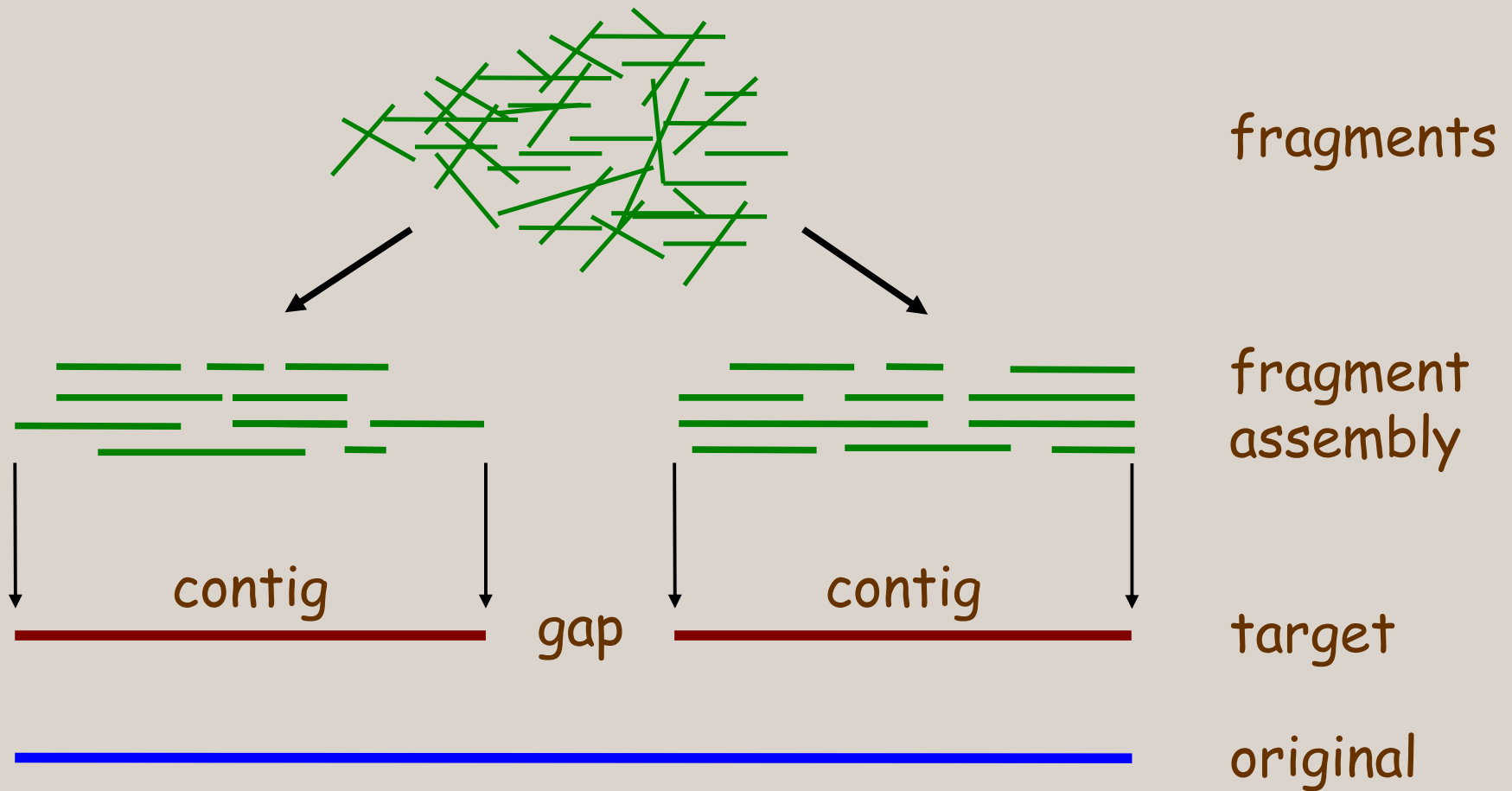Genomes can also be determined using a technique known as *shotgun sequencing*.

Computer scientists have played an important role in developing algorithms for assembling such data.

It's like putting together a jigsaw puzzle with millions of pieces (a lot of which are "blue sky").



Whole Genome Shotgun Sequencing Method

Genomic DNA

Sequence Each Fragment with Shotgun Approach

GCATTTCGAGTTACCTGGACAACCAGTG    GCTTGATTGGCCAATAATAGTATAT
CCAGTGGTACTGAGGACGCAAGAGGCTTGA

Align Contiguous Sequences

GCATTTCGAGTTACCTGGACAACCAGTGGTACTGAGGACGCAAGAGGCTTGATTGGCCAATAATAGTATAT

Generate Finished Sequence

http://occawlonline.pearsoned.com/bookbind/pubbooks/bc_mcampbell_genomics_1/medialib/method/shotgun.html

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Sequence Assembly



fragments

fragment assembly

contig          gap          contig

target

original
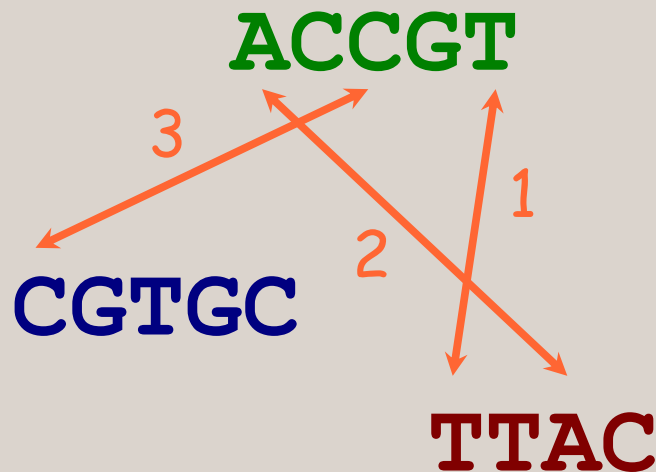
# Sequence Assembly

Simple model of DNA assembly is *Shortest Supersequence Problem*: given set of sequences, find shortest sequence $S$ such that each of original sequences is a subsequence of $S$.

Look for overlap between *prefix* of one sequence and *suffix* of another:

**ACCGT**

3

**CGTGC**

1

2

**TTAC**

```
--ACCGT--
----CGTGC
TTAC-----
_____
TTACCGTGC
```

# Sequence Assembly

Sketch of algorithm:

- Create an *overlap graph* in which every node represents a fragment and edges indicate overlap.

- Determine which overlaps will be used in final assembly: find an *optimal spanning forest* in overlap graph.
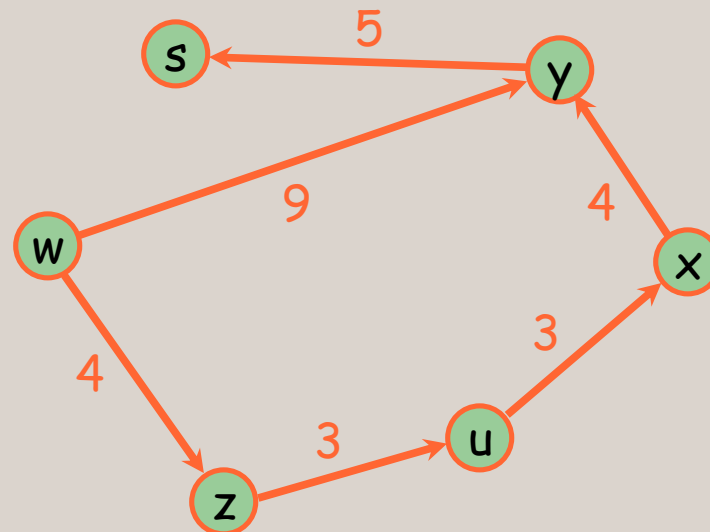
```
W = AGTATTGGCAATC

Z = AATCGATG

U = ATGCAAACCT

X = CCTTTTGG

Y = TTGGCAATCA

S = AATCAGG
```
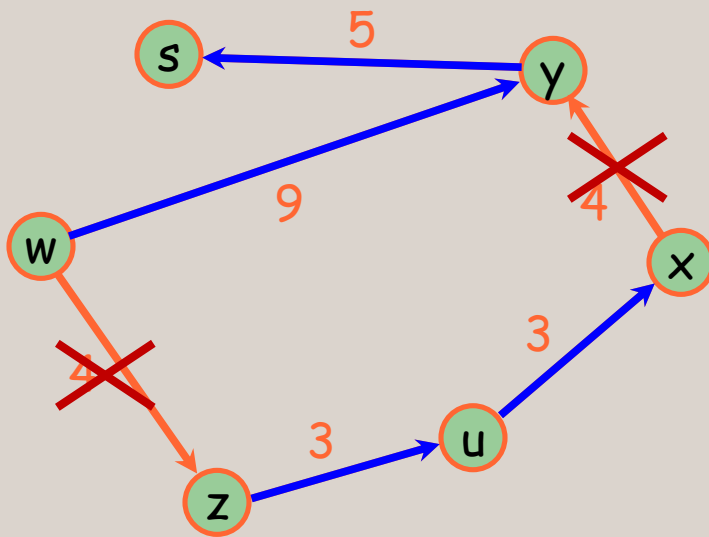
# Sequence Assembly

- Look for paths of maximum weight: use greedy algorithm to select edge with highest weight at each step.
- Edge must connect nodes with in- and out-degrees <= 1.
- May end up with set of paths: each yields a contig.



W→Y→S

```
AGTATTGGCAATC
    TTGGCAATCA
          AATCAGG
```
AGTATTGGCAATCAGG

Z→U→X

```
AATCGATG
    ATGCAAACCT
          CCTTTTGG
```
AATCGATGCAAACCTTTTGG

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Sequence Comparison

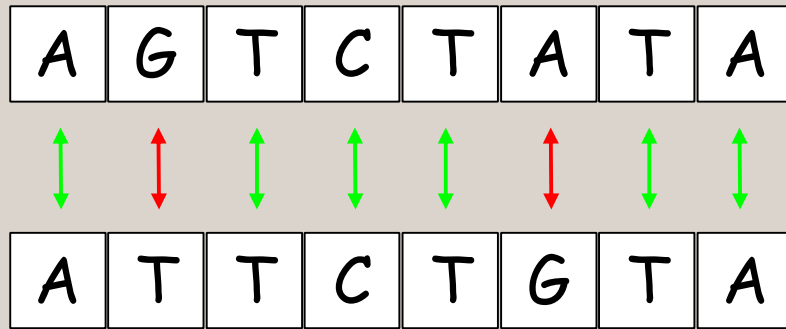What's the problem?  Kind of like google for biologists ...

- Given new DNA or protein sequence, biologist will want to search databases of known sequences for similarities.

- Sequence similarity can provide clues about function and evolutionary relationships.

- Databases such as GenBank are too big for manual search. To search them efficiently, we need an algorithm.

Can't expect exact matches (i.e., not really like google):

- Genomes aren't static:  mutations, insertions, deletions.

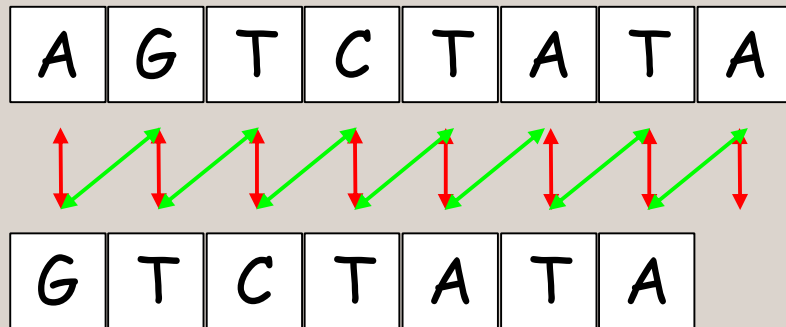- Human (and machine) error in reading sequencing gels.

HHMI
Howard Hughes Medical Institute

LEHIGH
U N I V E R S I T Y.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Sequence Comparison

Why not just line up sequences and count matches?

| A | G | T | C | T | A | T | A |
|---|---|---|---|---|---|---|---|

→ *Difference = 2*

| A | T | T | C | T | G | T | A |
|---|---|---|---|---|---|---|---|

Doesn't work well in case of deletions or insertions:

| A | G | T | C | T | A | T | A |
|---|---|---|---|---|---|---|---|

→ *Difference = 8*

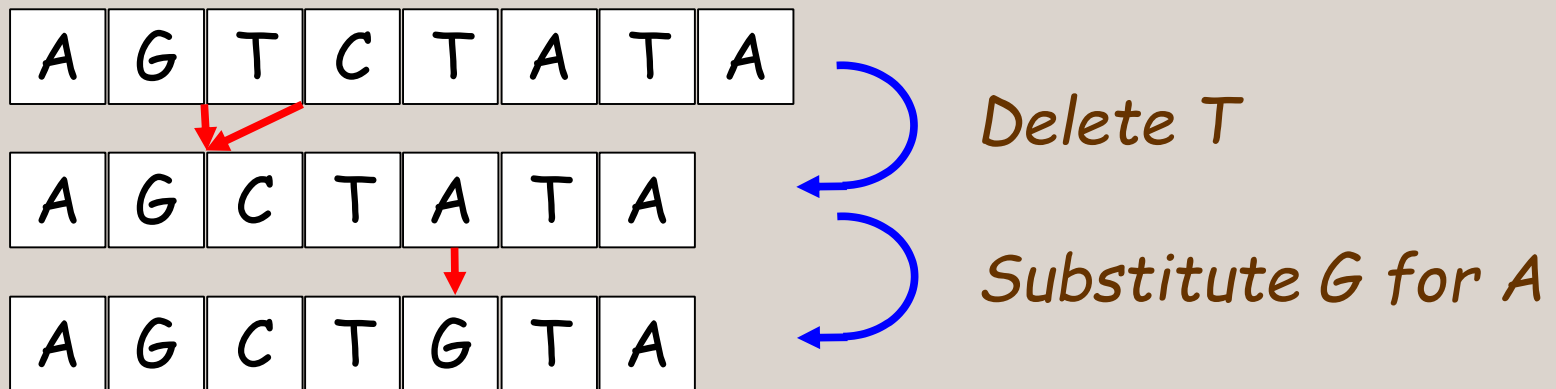| G | T | C | T | A | T | A |
|---|---|---|---|---|---|---|

One missing symbol at start leads to large difference!

# Sequence Comparison

Instead, we'll use technique known as *dynamic programming*.

- Three basic operations: delete a single symbol, insert a single symbol, substitute one symbol for another.

- Goal: given two sequences, find shortest series of operations needed to transform one into other.

| A | G | T | C | T | A | T | A |

*Delete T*

| A | G | C | T | A | T | A |

*Substitute G for A*

| A | G | C | T | G | T | A |

HHMI
Howard Hughes Medical Institute

LEHIGH
U N I V E R S I T Y.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
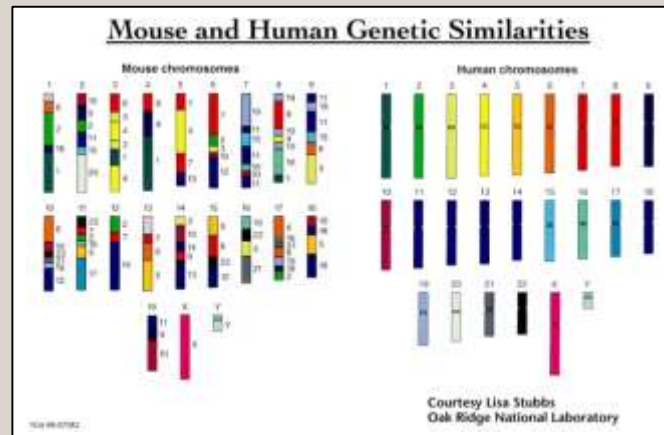Computer Science and Engineering
CSE

# Sequence Comparison

Elegant optimization algorithm builds table of values, working from shorter prefixes to longer prefixes:

ε ← → *s e q u e n c e   t*

ε

| 0 | ← cost of inserting *t* |

*sequence s* →

*cost of deleting s* ↑

$$d[i,j] = \min \begin{cases} d[i-1, j] + 1 \\ d[i, j-1] + 1 \\ d[i-1, j-1] + \begin{cases} 0 & \text{if } s[i] = t[j] \\ 1 & \text{if } s[i] \neq t[j] \end{cases} \end{cases}$$

HHMI
Howard Hughes Medical Institute

LEHIGH
U N I V E R S I T Y.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
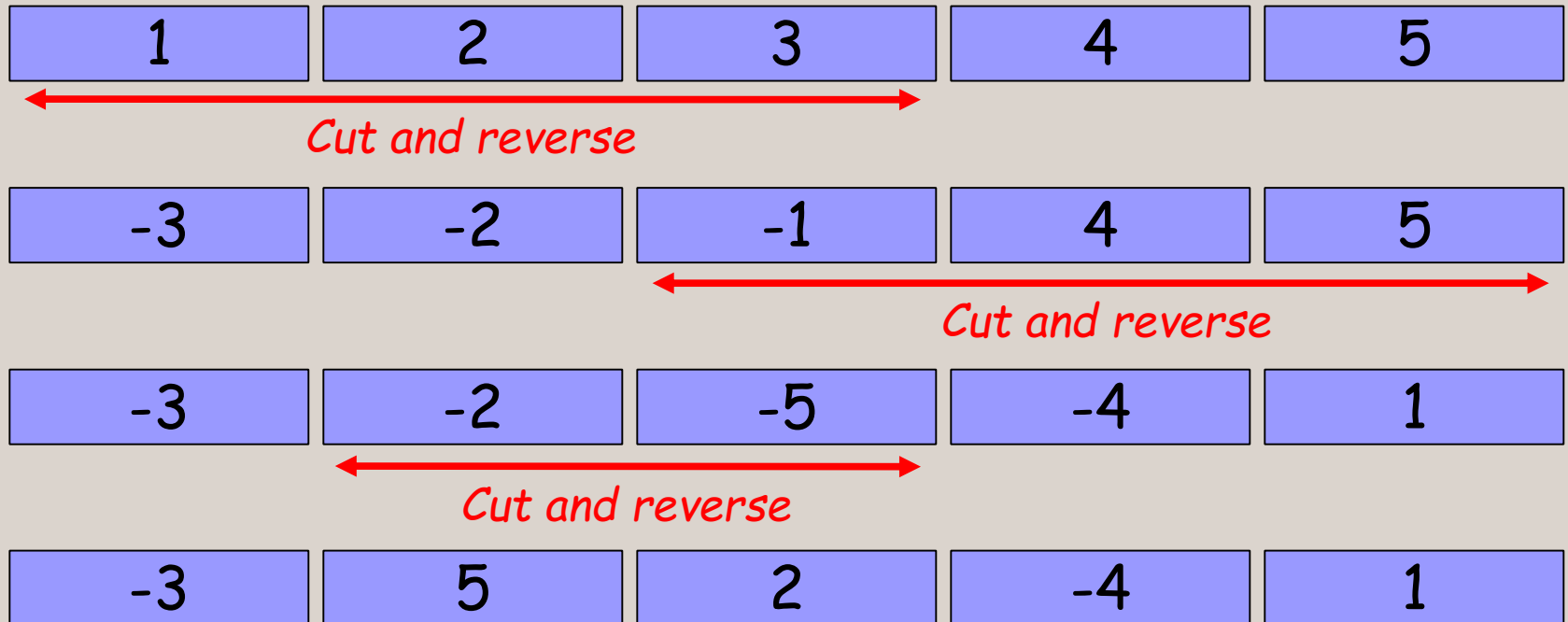CSE

# Genome Rearrangements

Recall what we saw earlier:



- 99% of mouse genes have homologues in human genome.

- 96% of mouse genes are in same relative location.

- Mouse genome can be broken up into 300 *synteny blocks* which, when rearranged, yield human genome.

- Provides a way to think about evolutionary relationships.

# Reversal Distance

Human Chromosome X

| 1 | 2 | 3 | 4 | 5 |

←————————————→
*Cut and reverse*

| -3 | -2 | -1 | 4 | 5 |

←————————————→
*Cut and reverse*

| -3 | -2 | -5 | -4 | 1 |

←————————————→
*Cut and reverse*

| -3 | 5 | 2 | -4 | 1 |

Mouse Chromosome X

*Reversal distance* is minimum number of steps needed.

# Interesting Sidenote

Early work on related problem, sorting by prefix reversals, was done in 1970's by Christos Papadimitriou, a professor now at UC Berkeley, and one "William H. Gates" ...



Yes, that Bill Gates ...

*Introduction to Bioinformatics* ▪ *Lopresti*
*BioS 10* ▪ *November 2019* ▪ *Slide 31*

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# History of Chromosome X



Hypothesized reversals

Rat Consortium, Nature, 2004

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

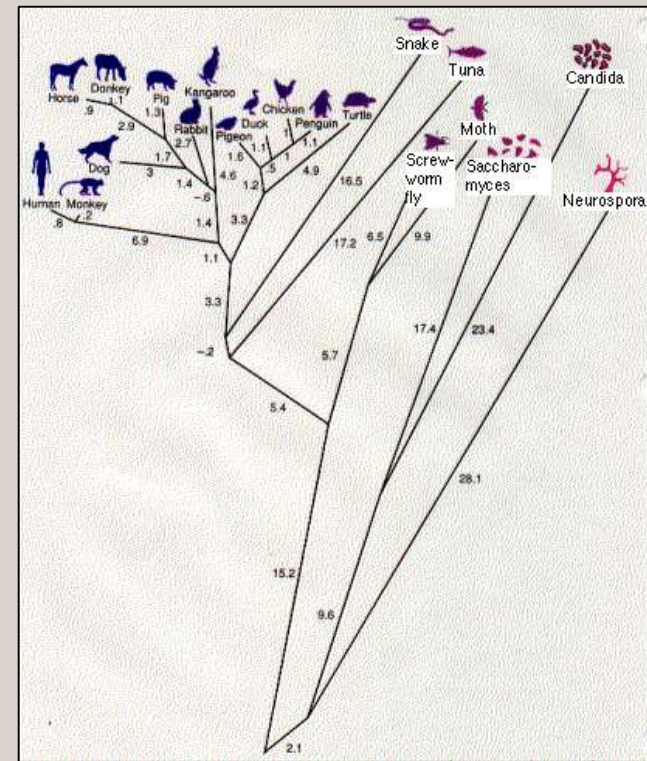# Building the "Tree of Life"

Scientists build phylogenetic trees to help understand evolutionary relationships.  Reversal distance often used.



Note:  trees are "best guesses" and certainly contain errors!

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# DNA Microarrays

- Allows simultaneous measurement of transcription level for every gene in a genome (gene expression).

- Differential expression, want to find genes that behave similarly over time.

- One microarray can test ~10k genes.

- Data obtained much faster than we can process it!

- Must find ways to uncover patterns.



*green = repressed*
*red = induced*

HHMI
Howard Hughes Medical Institute

LEHIGH
U N I V E R S I T Y.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Using DNA Microarrays



- Track sample over time to see change in gene expression.

- Track two different samples under same conditions to see difference in gene expressions.

*Each cell represents one gene's expression over time*

http://www.bioalgorithms.info/presentations/Ch10_Clustering.ppt

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# DNA Microarrays

*K-means clustering* is one way to organize this data:

- Given set of $n$ data points and an integer $k$.

- We want to find set of $k$ points that minimizes mean-squared distance from each data point to nearest center.

Sketch of algorithm:

- Choose $k$ initial center points randomly and cluster data.

- Calculate new centers for clusters using points in cluster.

- Re-cluster all data using new center points.

- Repeat second two steps until no data points change clusters, or some other convergence criterion is met.

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Clustering Microarray Data

- Pick *k* = 2 centers at random.
- Cluster data around these center points.

- Re-calculate centers based on current clusters.

From "Data Analysis Tools for DNA Microarrays" by Sorin Draghici.

HHMI
Howard Hughes Medical Institute

LEHIGH
U N I V E R S I T Y.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Clustering Microarray Data

- Re-cluster data around new center points.

- Repeat last two steps until no data points change clusters.

From "Data Analysis Tools for DNA Microarrays" by Sorin Draghici.

HHMI
Howard Hughes Medical Institute

LEHIGH
U N I V E R S I T Y.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Example of Hierarchical Clustering



*Different genes that express similarly*

From "Cluster analysis and display of genome-wide expression patterns" by Eisen, Spellman, Brown, and Botstein, Proc. Natl. Acad. Sci. USA, Vol. 95, pp. 14863–14868, December 1998

HHMI
Howard Hughes Medical Institute

LEHIGH
UNIVERSITY.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Why Study Bioinformatics?

- Many unanswered questions $\Rightarrow$ opportunities to make fundamental contributions (+ become rich and famous).

- Stretch your creativity and problem-solving skills.

- Cross-disciplinary teams: work with interesting people.

- Participate in unlocking the mysteries of life itself.

- Make the world a better place.

HHMI
Howard Hughes Medical Institute

LEHIGH
U N I V E R S I T Y.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Intro to Bioinformatics

Prof. Brian Chen

CSE 308 / BioE 308 covers:

- Intro to molecular biology & algorithms,
- Genetic sequence comparison & alignment,
- Sequencing & assembly of DNA,
- DNA microarrays,
- Gene regulatory networks,
- Genome annotation,
- Transcription factor binding site prediction,
- Standard formats and sources for genomic data, etc.

CSE 308 is <u>not</u> a programming course! It's for BioS, BioE, CSE, and Math students.

Questions:  chen@cse.lehigh.edu

HHMI
Howard Hughes Medical Institute

LEHIGH
U N I V E R S I T Y.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE

# Structural Bioinformatics

Prof. Brian Chen

CSE 307 / BioE 307 covers:

- Geometric modeling for proteins,
- Structure alignment & protein folding,
- Protein surfaces, cavities, electrostatics,
- Protein-protein and protein-DNA
- Interfaces and interactions,
- Protein structure prediction, simulation, docking,
- Structural bioinformatics in pharmaceutical discovery,
- Function annotation, active site prediction, etc.

*For seniors in BioS, BioE, CSE, and Math.*

Questions:  chen@cse.lehigh.edu

HHMI
Howard Hughes Medical Institute

LEHIGH
U N I V E R S I T Y.

Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
Computer Science and Engineering
CSE