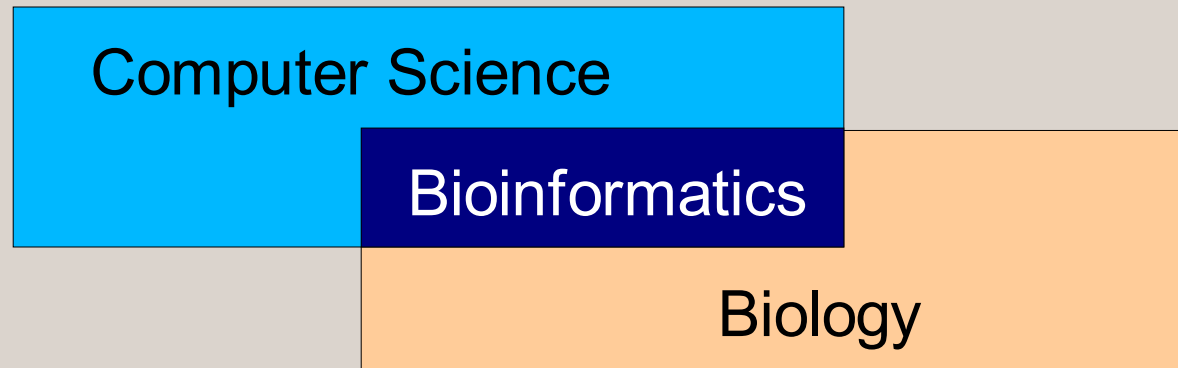# Life @ Lehigh

## Lehigh Life Days  •  April 12, 2006

*Dan Lopresti*

Associate Professor

Office PL 404B

dal9@lehigh.edu

LEHIGH
UNIVERSITY™

# Bioinformatics

What is bioinformatics?   Application of techniques from computer science to problems from biology.
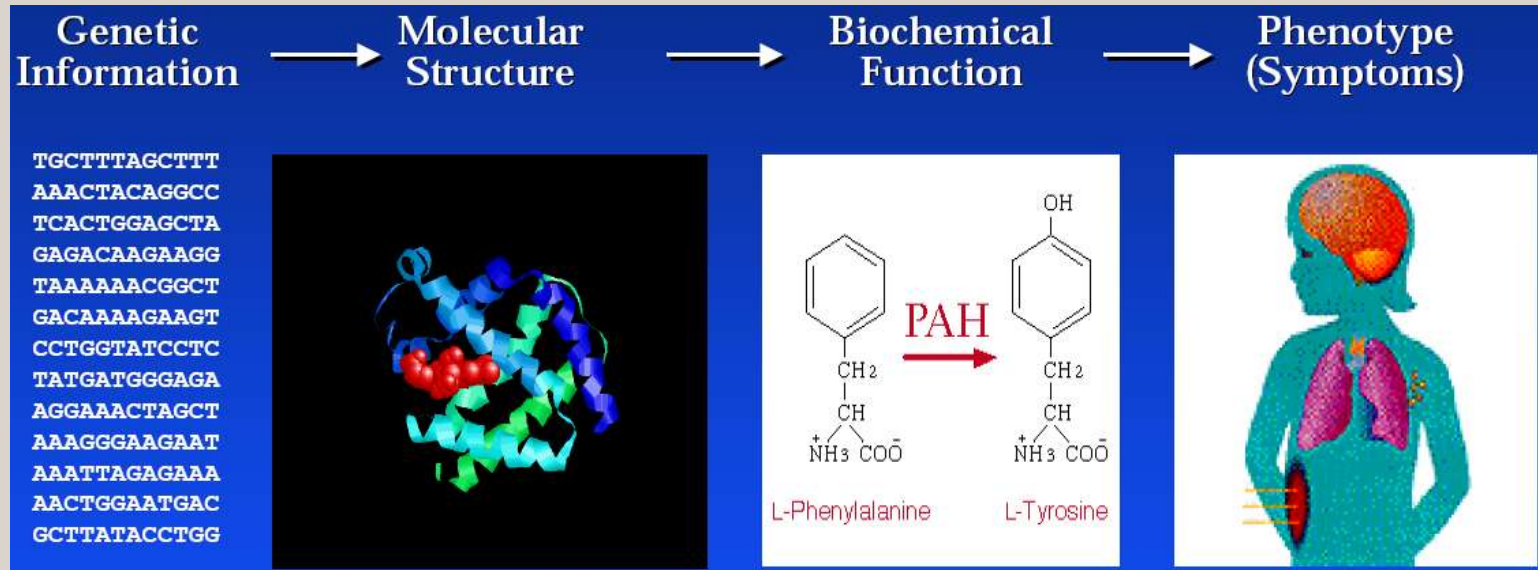
Computer Science

Bioinformatics

Biology

Why is it interesting?

- Important problems.
- Massive quantities of data.
- Desperate need for efficient solutions.
- Success is rewarded.

LEHIGH
U N I V E R S I T Y

# Motivation

"Biology easily has 500 years of exciting problems to work on."
*Donald Knuth*



By developing techniques for analyzing this data, we can attempt to understand genetic nature of diseases, evolution of life itself, etc.
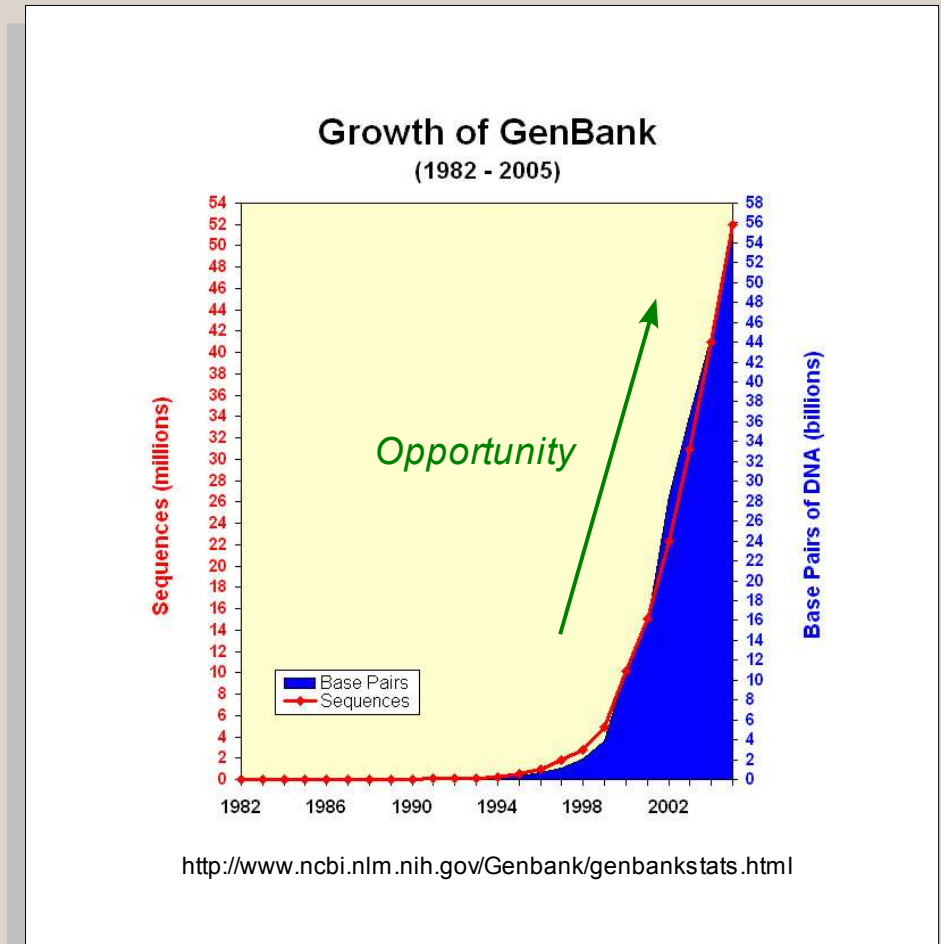
http://cmgm.stanford.edu/biochem218/

LEHIGH
UNIVERSITY

# Opportunity

Genetic identity of most organisms is encoded in long molecules made up of four basic units, the nucleic acids:

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

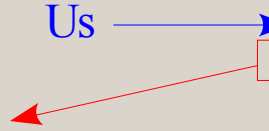To first approximation, DNA is language over 4 character alphabet, {A, C, G, T}.

**Growth of GenBank**
(1982 - 2005)

*Opportunity*

Base Pairs
Sequences

http://www.ncbi.nlm.nih.gov/Genbank/genbankstats.html

LEHIGH
U N I V E R S I T Y™

# Genomes

Complete 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 | 400,000,000 | 125,948,974 | 151,411 |
| Rattus norvegicus | 2,900,000,000 | 106,344,366 | 218,598 |
| Bos taurus | 3,651,500,000 | 71,215,626 | 159,473 |
| Glycine max | 1,115,000,000 | 62,817,102 | 141,802 |
| Medicago truncatula | 400,000,000 | 50,991,920 | 104,535 |
| Trypanosoma brucei | 35,000,000 | 49,855,996 | 91,334 |
| Lycopersicon esculentum | 655,000,000 | 49,415,566 | 97,112 |
| Giardia intestinalis | 12,000,000 | 47,639,714 | 54,328 |
| Strongylocentrotus purpur | 900,000,000 | 47,590,936 | 77,532 |
| Entamoeba histolytica | — | 44,522,016 | 49,938 |
| 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 |

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

LEHIGH UNIVERSITY

# Comparative Genomics



Mouse and Human Genetic Similarities
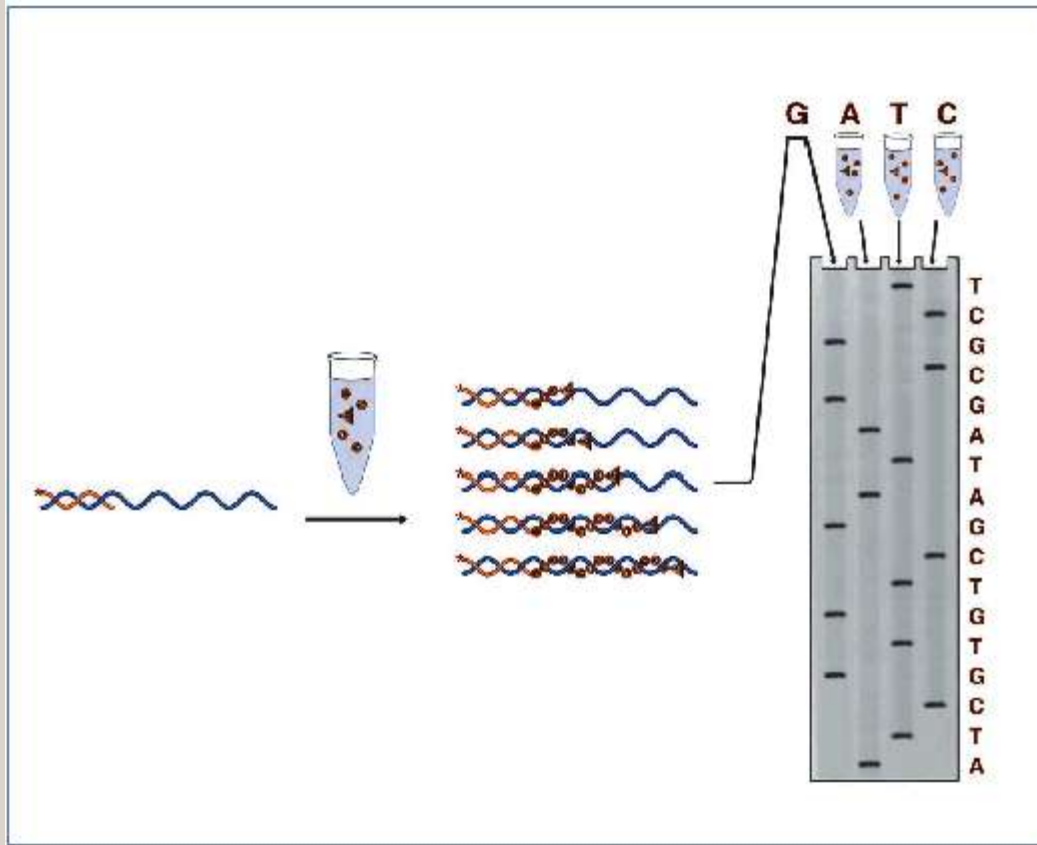
Courtesy Lisa Stubbs
Oak Ridge National Laboratory

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

LEHIGH
UNIVERSITY™

# Reading DNA



This is known as *Sanger sequencing.*

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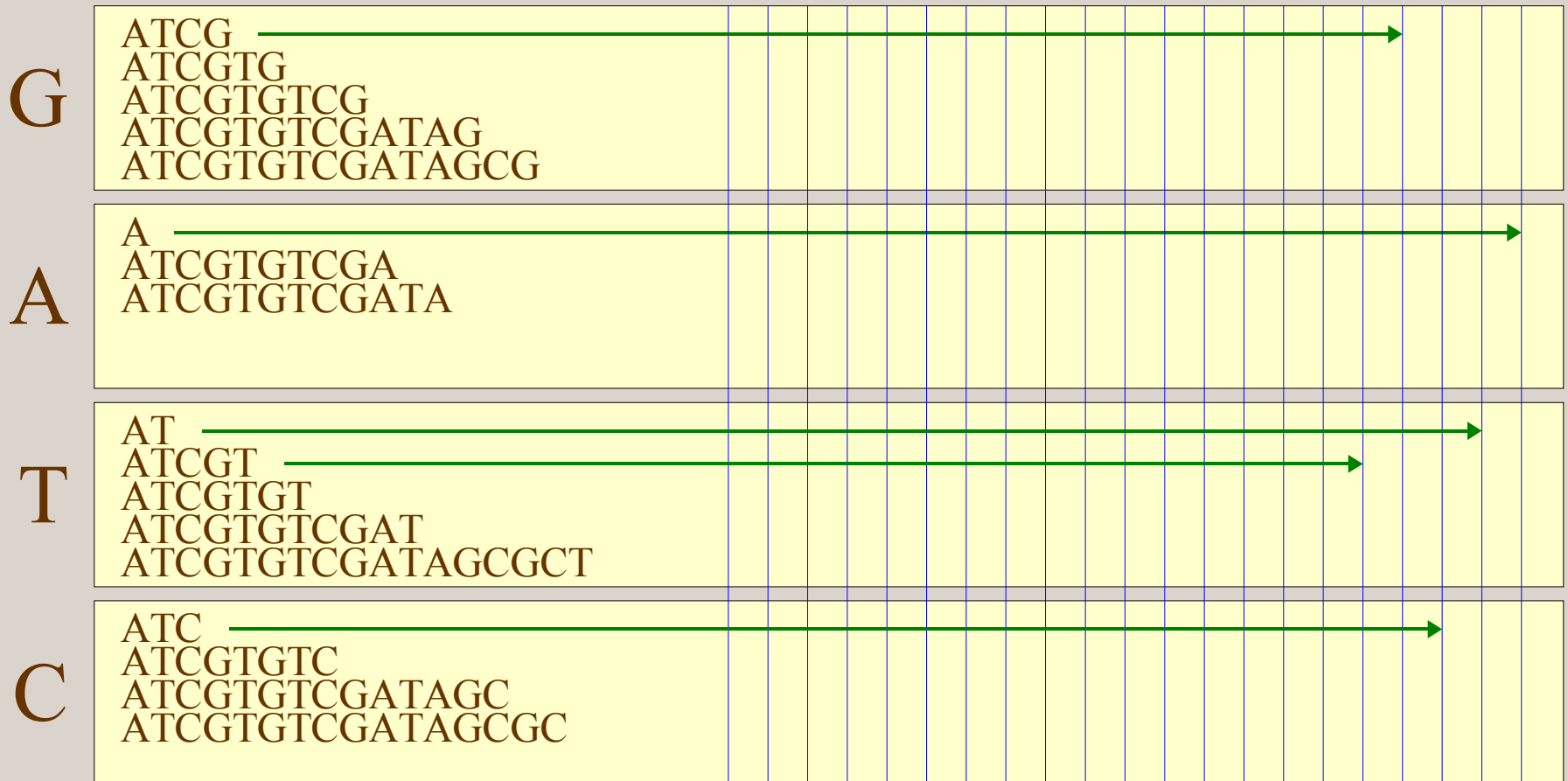http://www.iupui.edu/~wellsctr/MMIA/htm/animations.htm

*Gel electrophoresis* is process of separating a mixture of molecules in a gel media by application of an electric field.  In general, DNA molecules with similar lengths will migrate same distance.

First cut DNA at each base: A, C, G, T.  Then run gel and read off sequence: ATCGTG …

LEHIGH
U N I V E R S I T Y ™

# 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
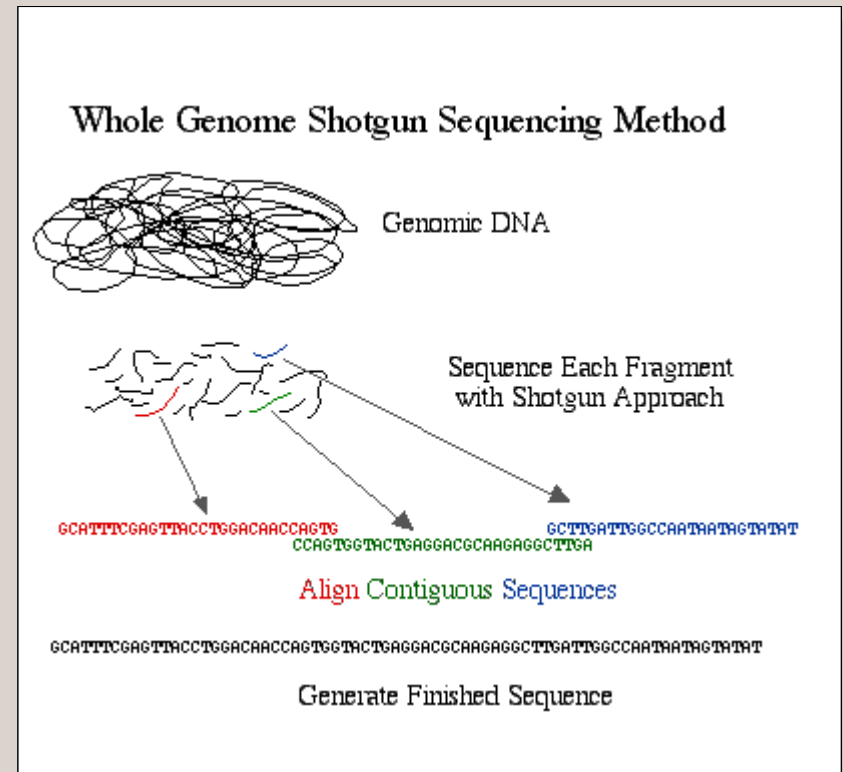ATCGTGTCGATAGCGC

LEHIGH
UNIVERSITY™

# Sequencing a Genome

Unfortunately, current sequencing technologies can only read 700 nucleotides at a time.
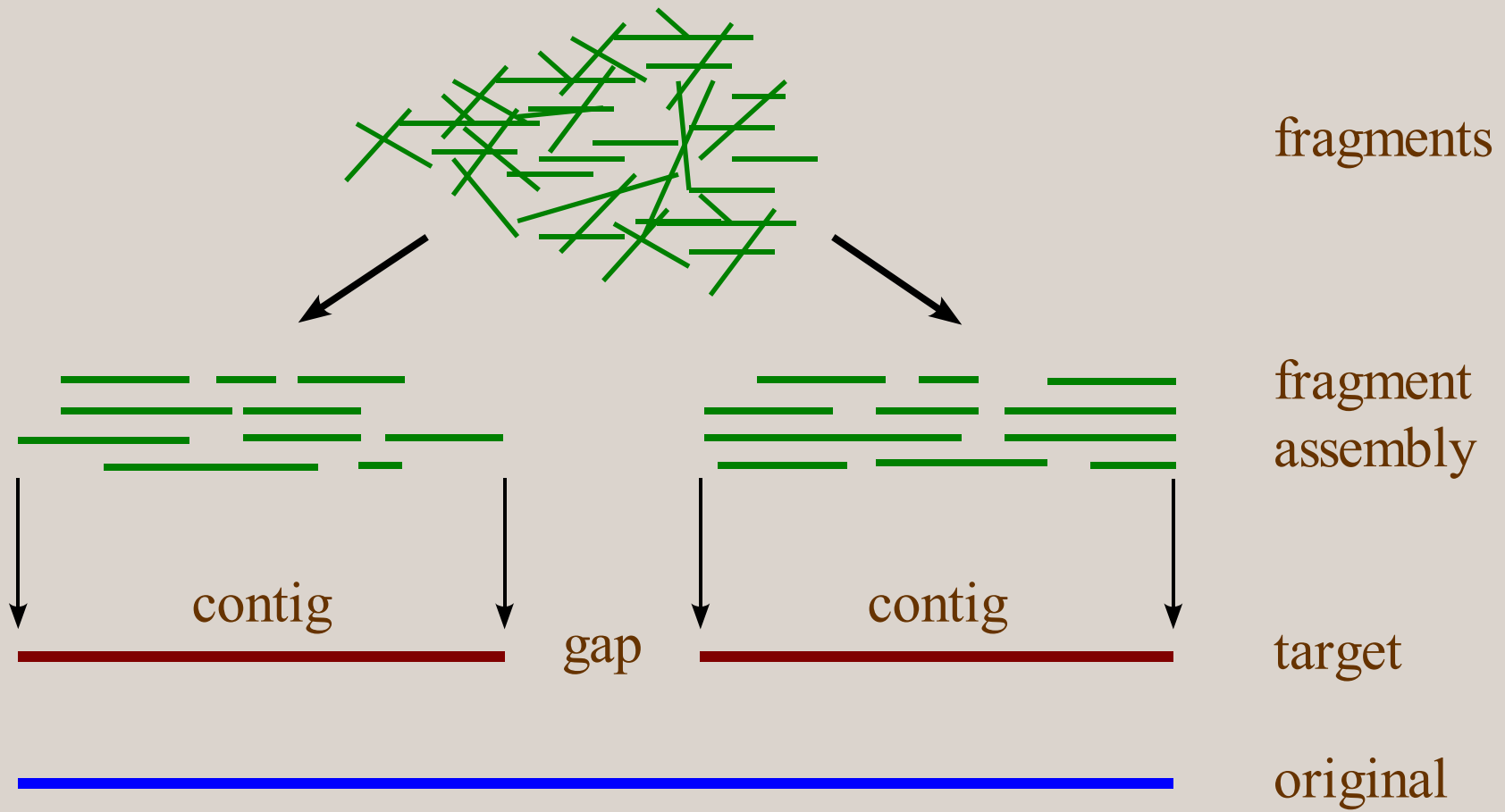
For genomes, we use *shotgun sequencing*, which breaks a chromosome into overlapping short sequences which must then be reassembled.

It's kind of 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
CCAGTGGTACTGAGGACGCAAGAGGCTTGA
GCTTGATTGGCCAATAATAGTATAT

Align Contiguous Sequences

GCATTTCGAGTTACCTGGACAACCAGTGGTACTGAGGACGCAAGAGGCTTGATTGGCCAATAATAGTATAT

Generate Finished Sequence

LEHIGH
UNIVERSITY

# Sequence Assembly



fragments

fragment assembly

contig

contig

gap

target

original

LEHIGH
UNIVERSITY™
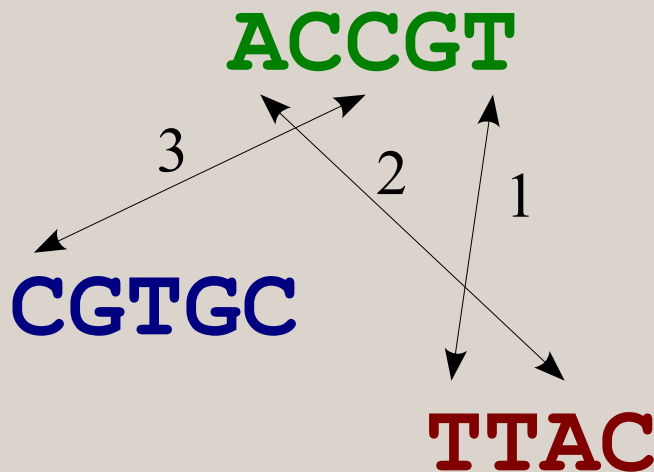
# Sequence Assembly

A simple model of DNA assembly is the *Shortest Supersequence Problem*:  given a set of sequences, find the shortest sequence *S* such that each of original sequences appears as subsequence of *S*.

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

**ACCGT**

3      2      1

**CGTGC**

**TTAC**

```
--ACCGT--

----CGTGC

TTAC-----
_____
TTACCGTGC
```

LEHIGH
UNIVERSITY

# Sequence Assembly

Sketch of algorithm (procedure for assembling fragments):

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

● Determine which overlaps will be used in the final assembly: find an optimal collection of paths in overlap graph.
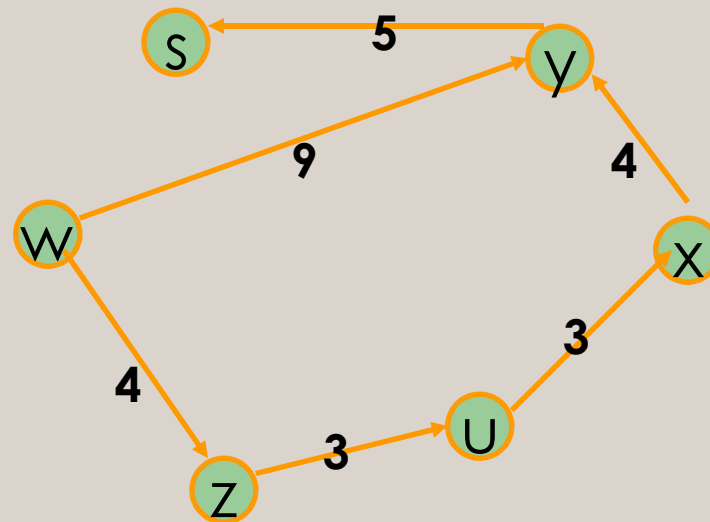
```
W = AGTATTGGCAATC

Z = AATCGATG

U = ATGCAAACCT

X = CCTTTTGG

Y = TTGGCAATCA

S = AATCAGG
```
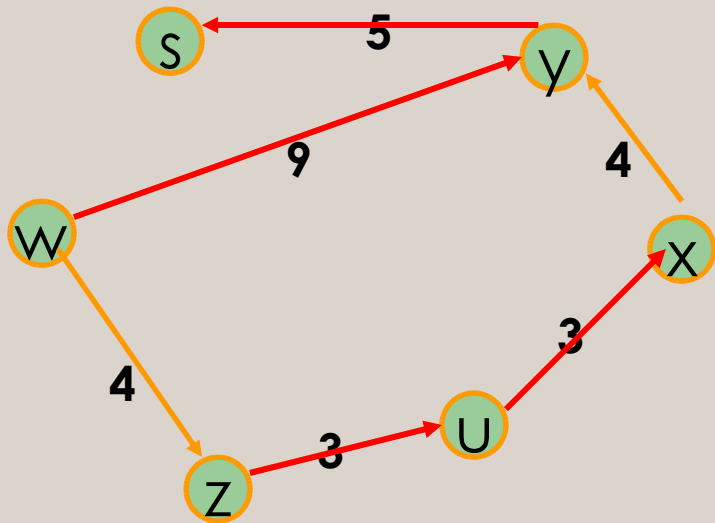
# Sequence Assembly

- Look for paths of maximum weight:  use *greedy* algorithm to select edge with highest weight at every step.

- Selected edge must connect nodes with in- and out-degrees <= 1.

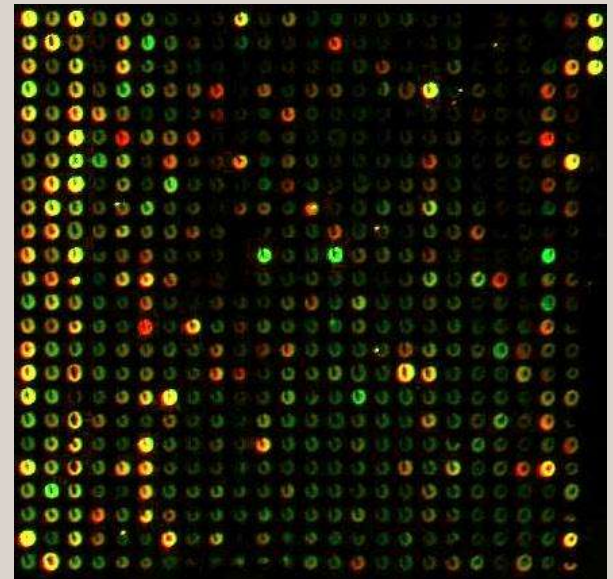- May end up with set of paths:  each corresponds to a contig.



$W{\rightarrow}Y{\rightarrow}S$

```
    AGTATTGGCAATC
        TTGGCAATCA
              AATCAGG
```
AGTATTGGCAATCAGG

$Z{\rightarrow}U{\rightarrow}X$

```
AATCGATG
    ATGCAAACCT
          CCTTTTGG
```
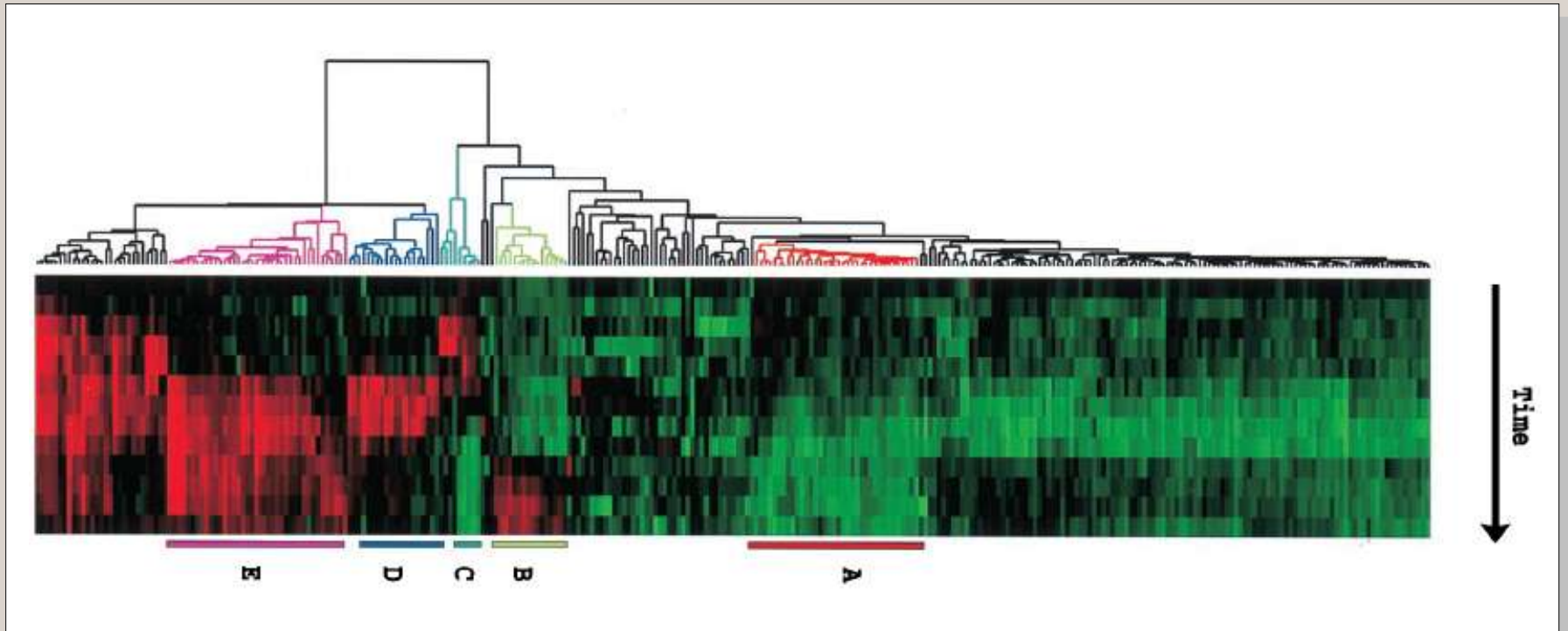AATCGATGCAAACCTTTTGG

LEHIGH
UNIVERSITY

# DNA Microarrays

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

- Differential expression, changes over time.

- Single microarray can test ~10k genes.

- Data obtained faster than can be processed.

- Want to find genes that behave similarly.

- A pattern discovery problem.

*green = repressed*

*red = induced*

LEHIGH
U N I V E R S I T Y ™

# Visualizing Microarray Data



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

LEHIGH
UNIVERSITY™

# Clustering Microarray Data

*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 $k$ center points that minimizes mean-squared distance from each data point to its nearest cluster center.
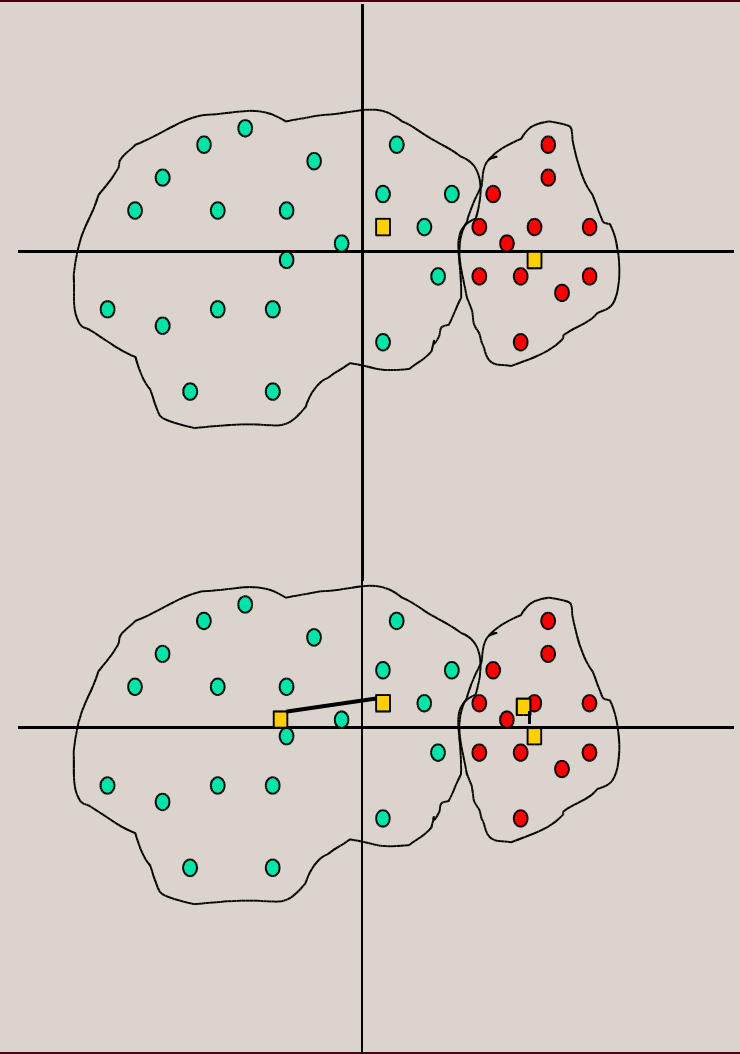
Sketch of algorithm:

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

- Calculate new centers for each cluster using points in cluster.

- Re-cluster all data using new center points.

- Repeat second two steps until no data points are moved from one cluster to another or some other convergence criterion is met.

LEHIGH
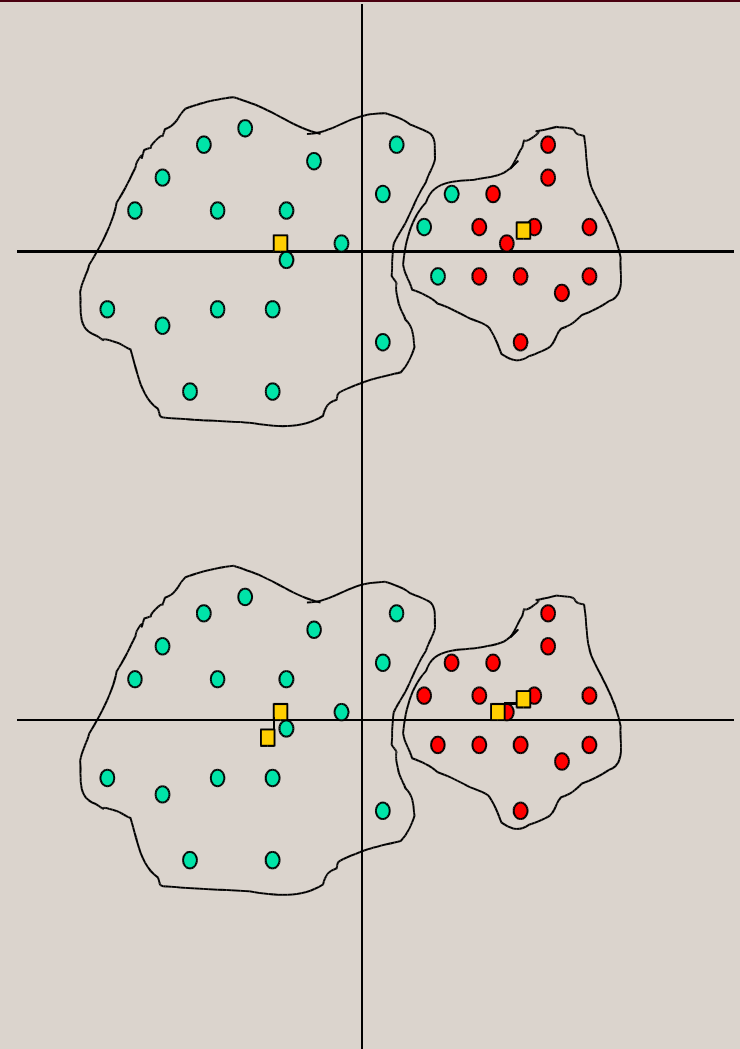UNIVERSITY™

# 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.

LEHIGH
U N I V E R S I T Y ™
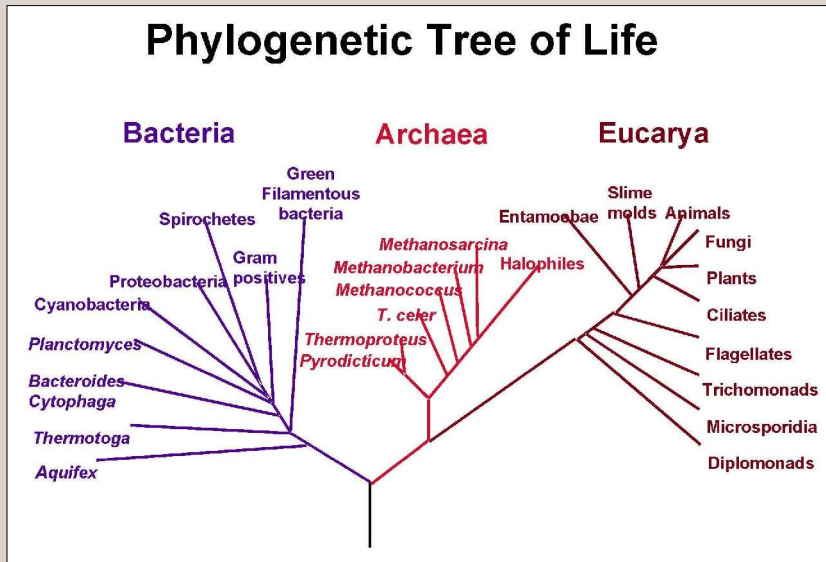
# Clustering Microarray Data

- Re-cluster data around new center points.

- Repeat last two steps until no more data points are moved into a different cluster.
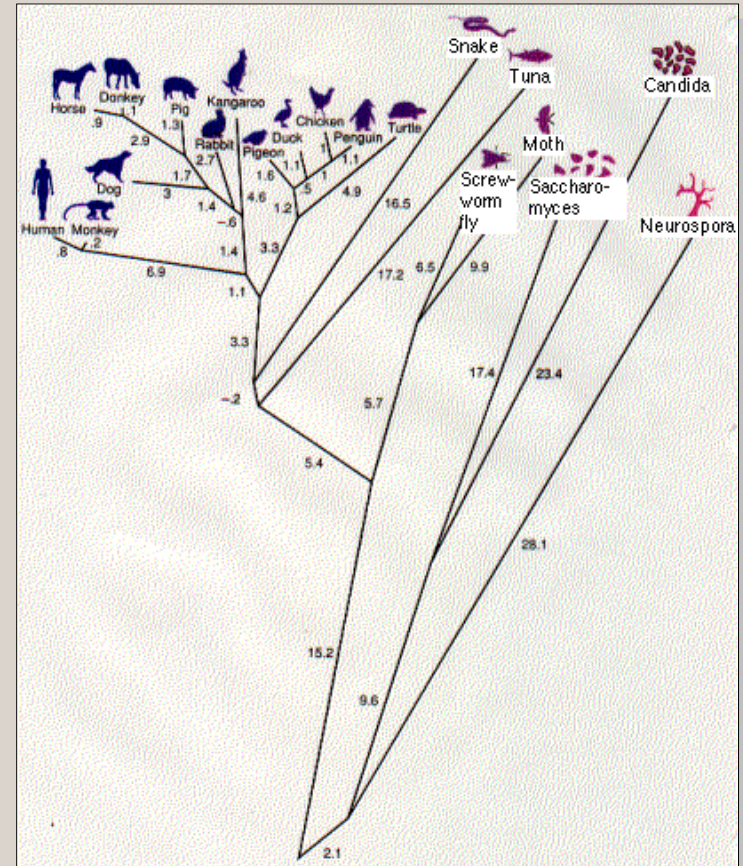
From "Data Analysis Tools for DNA Microarrays" by Sorin Draghici.

LEHIGH
U N I V E R S I T Y ™

# Building the "Tree of Life"

Scientists build phylogenetic trees in an attempt to understand evolutionary relationships.



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

LEHIGH
UNIVERSITY

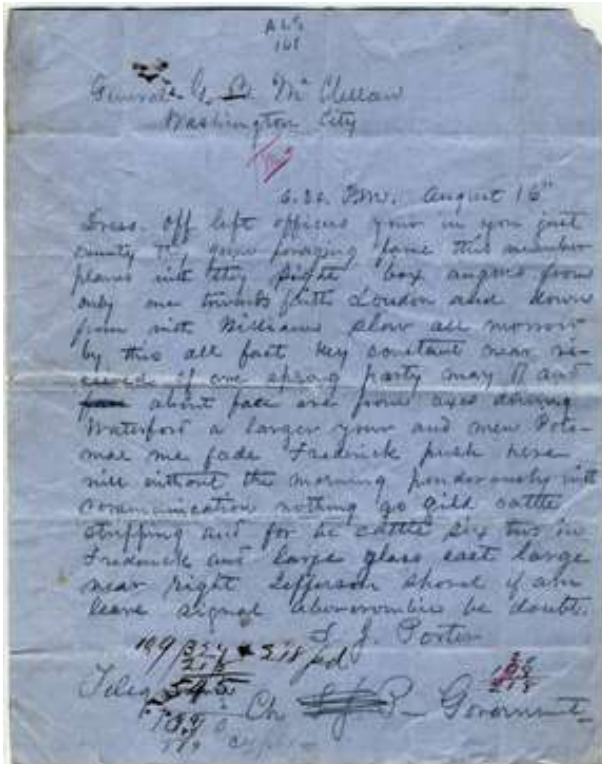# Why Study Bioinformatics?

- Still many urgent open problems $\Rightarrow$ lots of opportunities to make fundamental contributions (and become rich and famous).

- Stretch your creativity and problem-solving skills to the limit.

- Join a cross-disciplinary team – work with interesting people.

- Participate in unlocking the mysteries of life itself.

- Make the world a better place.

LEHIGH
UNIVERSITY

# My Advice on Colleges

- Gather all the data. Then trust your instincts.

- Whatever your decision, college will probably be the best years of your life – take full advantage of it.

- Beyond the obvious criteria, look for opportunities to become engaged in cutting-edge research, whatever your major.

- Seek out projects that will attract attention and have an impact.

LEHIGH
UNIVERSITY™

# Breaking a Civil War Secret Code

## The Civil War letter ...



http://digital.lib.lehigh.edu/remain

... encrypted – not yet broken.

## The players ...



**Major General
Fitz-John Porter**
author of letter, blamed for
Union loss at Second Bull Run,
court-martialed in 1863

**Major General
George McClellan**
recipient of letter



**Anson Stager**
inventor of cypher system
and later an early leader in
U.S. telecommunications

LEHIGH
U N I V E R S I T Y

# Breaking a Civil War Secret Code

## The news story ...



"Lehigh team works to crack Civil War code"

... finale yet to be written.
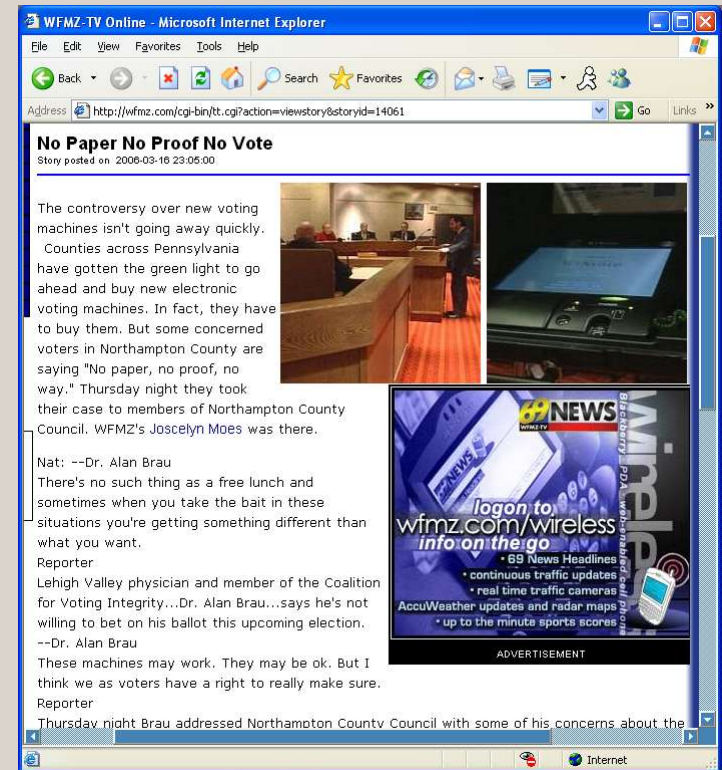
## The current software tool ...



... promising, but needs work.

Hint:  I'm looking for help ...

LEHIGH
UNIVERSITY

# Better Electronic Voting

E-voting has generated some controversy recently ...



Md. House Approves Paper Ballots - Microsoft Internet Explorer

**Md. House Approves Paper Ballots**

By Ann E. Marimow
Washington Post Staff Writer
Friday, March 10, 2006; Page B04

The Maryland House of Delegates unanimously passed legislation yesterday to ditch the state's touch-screen voting machines for the coming election in favor of a system that uses paper ballots.

The 137 to 0 vote in the House and the endorsement of the plan this week by Republican Gov. Robert L. Ehrlich Jr. represents a stunning turnaround for a state that was on the leading edge of touch-screen voting in 2001, and it reflects a national shift toward machines that provide a paper record.

The touch-screen system, for which Maryland has committed more than $90 million, would be put aside for one year while the state spends at least $13 million to lease optical scan machines.

"It's critically important for voters to know their vote was cast and that it will be counted correctly," said Del. Obie Patterson (D-Prince George's).

Linda H. Lamone, Maryland's top elections official, says paper ballots provide a "false



WFMZ-TV Online - Microsoft Internet Explorer

**No Paper No Proof No Vote**
Story posted on 2006-03-16 23:05:00

The controversy over new voting machines isn't going away quickly. Counties across Pennsylvania have gotten the green light to go ahead and buy new electronic voting machines. In fact, they have to buy them. But some concerned voters in Northampton County are saying "No paper, no proof, no way." Thursday night they took their case to members of Northampton County Council. WFMZ's Joscelyn Moes was there.

Nat: --Dr. Alan Brau
There's no such thing as a free lunch and sometimes when you take the bait in these situations you're getting something different than what you want.
Reporter
Lehigh Valley physician and member of the Coalition for Voting Integrity...Dr. Alan Brau...says he's not willing to bet on his ballot this upcoming election.
--Dr. Alan Brau
These machines may work. They may be ok. But I think we as voters have a right to really make sure.
Reporter
Thursday night Brau addressed Northampton County Council with some of his concerns about the

Maryland votes "yes" for paper trail ...     while Pennsylvania votes "no."

LEHIGH UNIVERSITY

# Better Electronic Voting
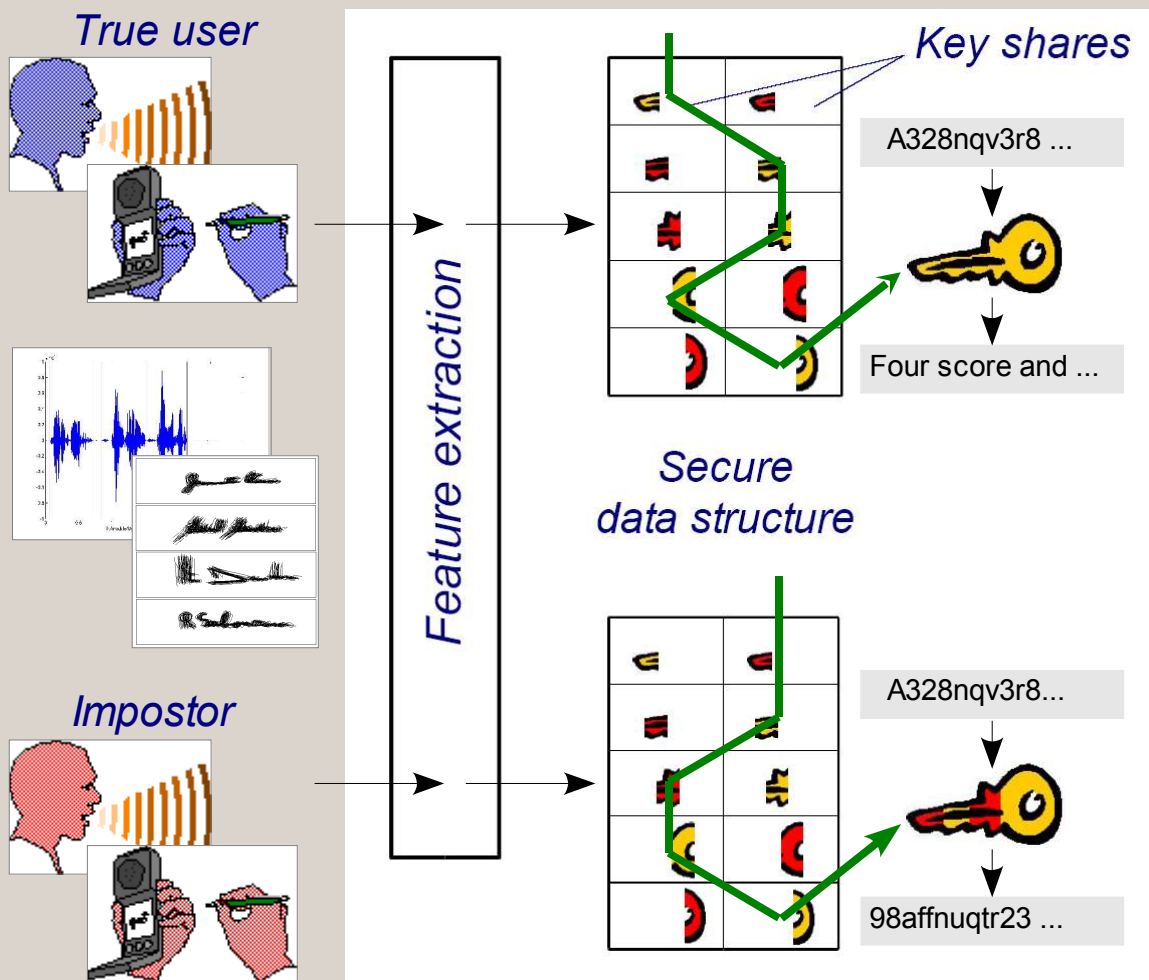
E-voting:  what's the right answer?

- Take a critical look at all aspects of the problem.

- Examine both security and usability issues.

- Build a prototype of an e-voting system that includes a reliable Voter Verified Paper Audit Trail (VVPAT).

- Some critics claim it can't be done: we disagree.

- Yet another undergraduate project.

- Of fundamental importance because our demoncracy depends on fair and transparent elections.



Diebold e-voting system

LEHIGH
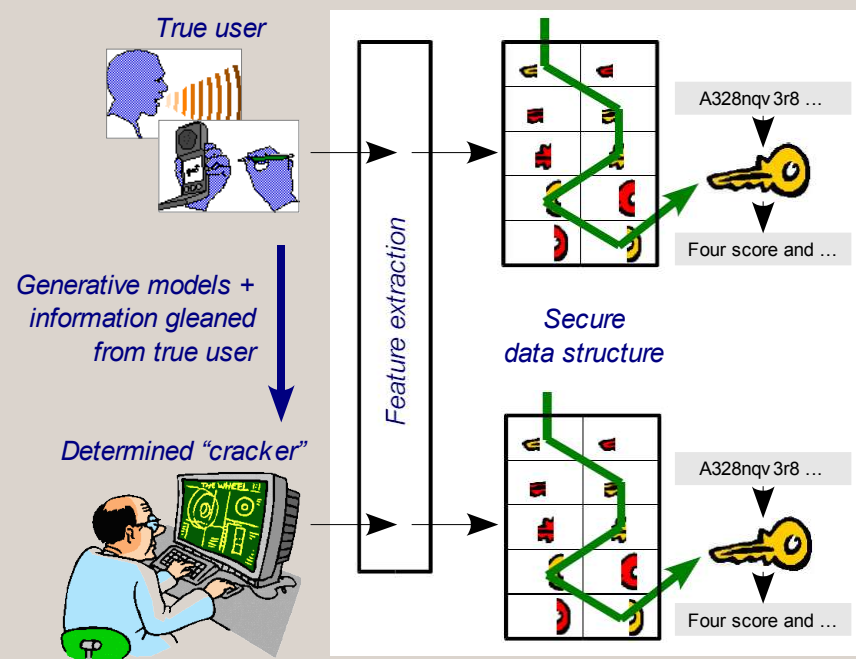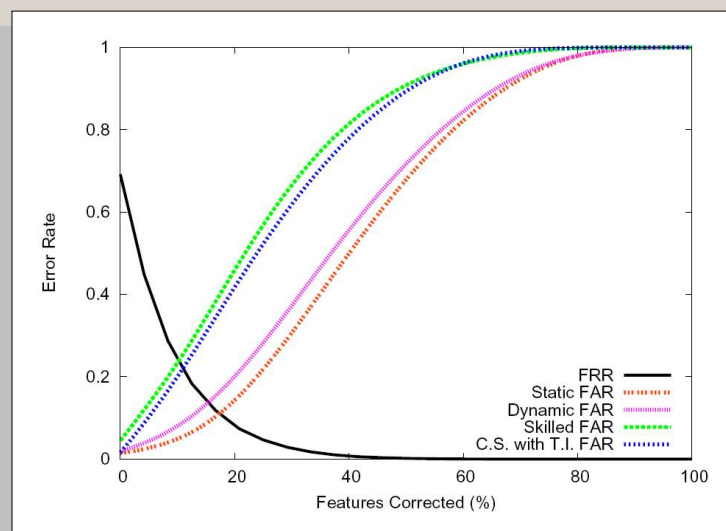UNIVERSITY™

# Evaluating Biometric Security

- Cryptographic key broken into shares and mixed with random data.

- Features extracted from user's speech or handwriting.

- Only input from true user will select correct shares to yield proper key.



True user

Key shares

A328nqv3r8 ...

Four score and ...

Feature extraction

Secure data structure

Impostor

A328nqv3r8...

98affnuqtr23 ...

LEHIGH
UNIVERSITY™

# Evaluating Biometric Security

Biometrics may be vulnerable to attacks using generative models.

- Some current systems at risk.
- Results for handwriting show machine can equal performance of skilled human forger:





Use our experience to improve biometrics, increase security.

LEHIGH
UNIVERSITY

# Last but not least ...

Why attend Lehigh?

Because this is a great place to be a student!

# Thank you!

LEHIGH
UNIVERSITY™