# CMSC423:
# Bioinformatic databases, algorithms and tools

*Héctor Corrada Bravo*
Dept. of Computer Science
Center for Bioinformatics and Computational Biology
University of Maryland

*University of Maryland, Spring 2017*

# Advances in Biology and Medicine needed, need, and will continue to need computational and statistical thinking
# (and their tools)

*Héctor Corrada Bravo*
Dept. of Computer Science
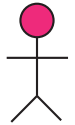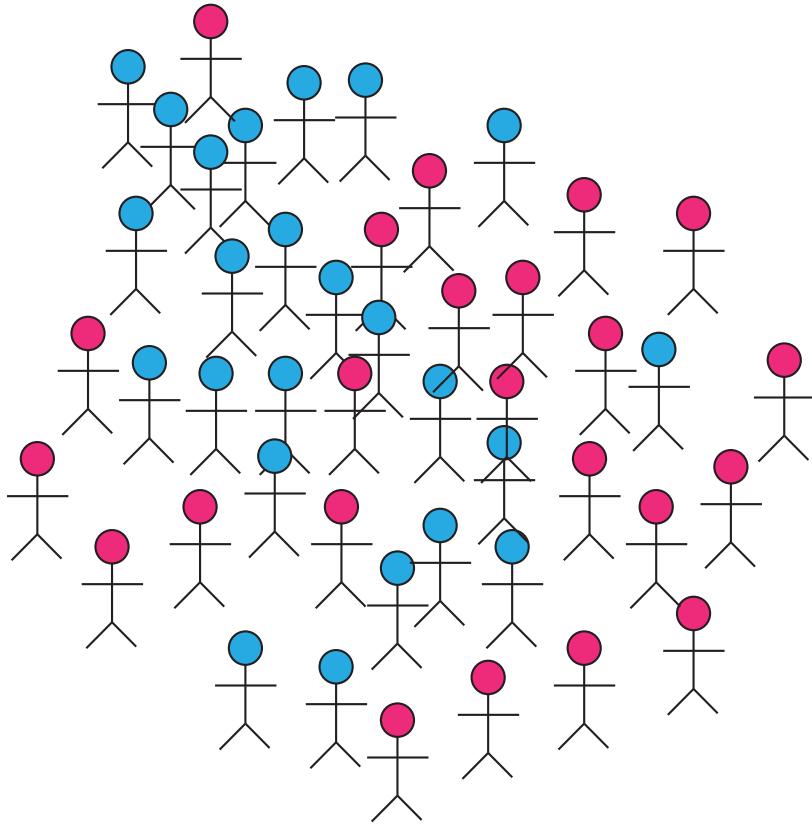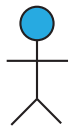Center for Bioinformatics and Computational Biology
University of Maryland

# What is Genomics?

- Each cell contains a complete copy of an organism's genome, or blueprint for all cellular structures and activities.

- The genome is distributed along chromosomes, which are made of compressed and entwined DNA.

- Cells are of many different types (e.g. blood, skin, nerve cells), but all can be traced back to a single cell, the fertilized egg.

# What is Genomics?

- Study the **molecular** basis of *variation* in development and disease

- Using **high-throughput** experimental methods

  - algorithms

  - ML

  - data management
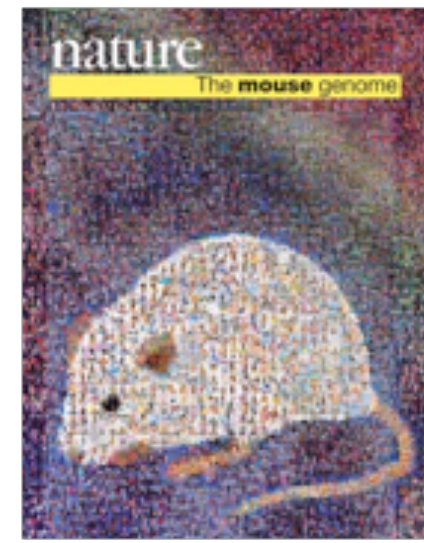
  - modeling

cancer

healthy

# Measurement

- For a small enough piece, we can measure the sequence of bases, referred to as *sequencing*
- Human Genome Project



D. melanogaster, Science, 2000

H. sapiens, Nature, 2000
and Science, 2000

M. musculus, Nature, 2002

# Genome

```
TCAGTTGGGAGCTGCTCCCCCACGGCCTCTCCTCACATTCCACGTCCTGTAGCTCTATGACCTCCACCTTTGAGTCCCTCCTCTCACACCTGAC
ATGAAAAGGCACATGAGGATCCTCAAATACCCCGTGATCAGTCTCAGGGTAGCTCTCATAGCCTGGACAGGGCCCCCCTCGGGGGTTGCGCCC
AGGTCCAGGCGGGGGATGCACAGCAACAGTCACCGAAGCAGAAGCCGTCACAGTGGTGATGGGCTGGCAGTAGCTGGGCACAGAGCTGCCCAT
GGCGGTGGACGTTGGGTTCCGAGGGTTGTGAGAACGGGCCCCACGGGGCCCTGAGCGGTCCCTATTGCTAGGGCCAGAATGCCCTTCAGTAGA
AATTTCAAAAGCGTCTCTGCGCGGTCTGTAGGGGGGTGGCCGCAAGCCTTCTCTAGGGGGATCCCTTCGAGGCTGCTGGCCTTGCCGTCCAGG
GGACAAGGAGCCAGAGTCCAGGTGGGGCTGTTGCCGAGGGGTCAAGGGAGGCTGATGTCTGGAGTCCGGATGGACCACCTGCAGAGGAGAGAC
ATAGGTCAACACAGGGAGGTAGGATGGTGGTGATGTTCCACCCACAAAAGAAAACCTATTCCTTTAGAAACCTCCAGGATGTGAATCCTGCCT
GCACCTGCACAGCTGGCTGGAGGCATATAGCCACTGCCCATAGATCTCAACTTACCCTCACAACCAACTGCCCCCAGGCCTAAGTTCTCTGCC
TCAAAACTGCCAAGGCCTGGATAGCCAAGAGCCTGGGTGTCTTGGAAATATGCAACCATAAATAGTAGCTTTTAGAAGTATAAGGCTCCTGTT
TCTGGGTCATATTAGTGTTGTTTTCACCTGTCCCCAGCCCTAAGCCAGGTGTGGCCAGAAGCAAATGTACTGTAAGAGCAGAGCAAAAACTTC
CACACAGATAGTTCTGTTAGGCAATACATCTCTGCCTGACTATTAGGAATCTGGTTTCTGGGTCCTCTGTACAAAGCTCGGAGCAACACAGTG
GCCACATCAATCAAAAGGACCGTGACCAACTTCAAAGTCGGTGAGCTTGTACCTATTTTTAGGCTCCTGCTGAACAGAACCAGATTCACACTA
CAGCTCAGCAGGGCATCGTCACGGGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTGTTGGGGGGGGGGGGTGGACAGAGGACGGGGAC
ACAATTCACTGGCCAGCCCTTCTCTCCTTCAAGGAAGGCTGCTCTAGCCTGGGACTGGAATACACATTTCCTGTAAACATGGTGGGGGCCTCA
GGCAAGCCAGAGTTTTGGAGCCTTCCTTAACTCTTCAAGGTGAGCATCTTGACTTGGAGGGTGGGGGTGCGGGTAAGGAAGGAACCTGTGGAC
TCCTCCCTACAAGACAGAAAAGGAATAAGCCACGAAGACAATAACGATTTTTGTATCAAGCGTCCTCTCCCATTTCAGCTTACCTGACAATGA
AATCAAATTCGGACCCTGCAAGCATCAGTACACCCAGCAGAGTGGACACAGCACCGTCCAGAACGGGAGCAAACATGTGCTCCAGAGCGAGCA
TAGCCCTGTGGTTCTTGTCCCCAATGGCTGTCAGAAAGGCCTGAACAAAGGAGAAATTGACACGGTCACATTCTGGGTGTGGTAAAGTGCTC
AGCTGTGTCTATACTTGGGTTTTGTAT...
```

**Total amount of DNA in human genome:**
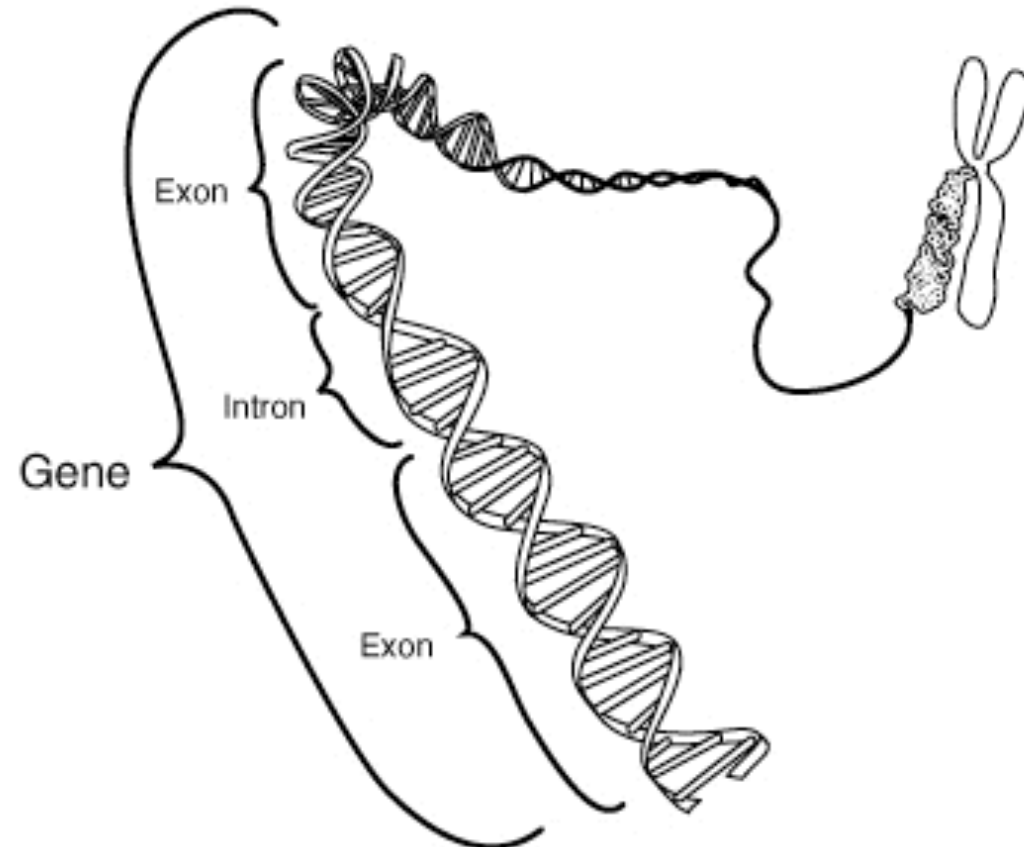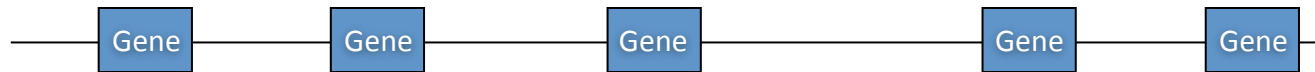**3 * $10^9$ base pairs (bp)**

# Why are these two different?



Differences explained by 1-10% difference in genome

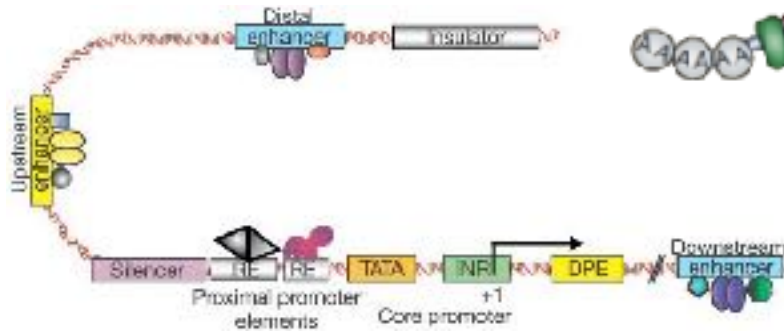Similarities explained by similar genes
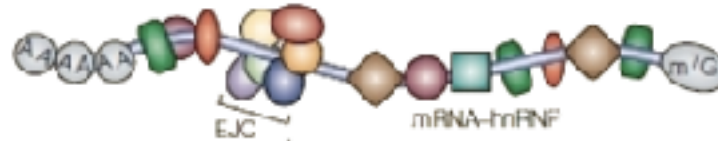
# Genes

# Computational Biology

Genes encode proteins which are transcribed into mRNA and translated into proteins.
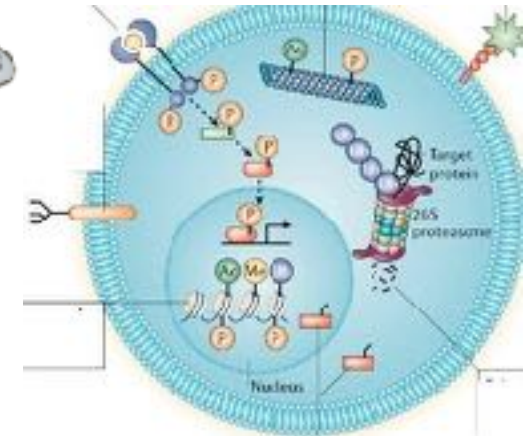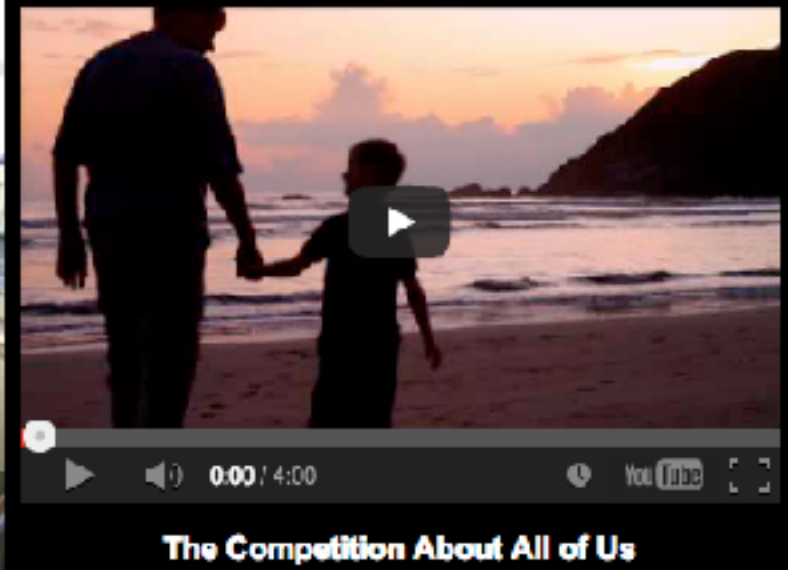
genomics                    transcriptomics                    proteomics



**Major** technological advances allow **unprecedented** data acquisition

ARCHON GENOMICS XPRIZE

PRESENTED BY
EXPRESS SCRIPTS

Follow us

SUBSCRIBE

COMPETITION DETAILS    XPRIZE 100 OVER 100 CANDIDATES    TEAMS    MEDIA    BLOG    ABOUT    SCIENTIFIC COMMUNITY

100 Over 100 Candidate Irving Kahn, Featured on CNN

0:00 / 4:00

The Competition About All of Us

build a **whole human genome** sequencing device and use it to sequence **100 human genomes** within **30 days or less**, with an accuracy of no more than one error in every 1,000,000 bases sequenced, with an accuracy rate of at least 98% of the genome, and at a recurring cost of **no more than $1,000 (US) per genome**.

NBC NEWS SCIENCE

COSMIC LOG
$10 million Genomics X Prize canceled; 'Outpaced by innovation'

# $10 million Genomics X Prize canceled: 'Outpaced by innovation'

Alan Boyle, Science Editor, NBC News

Aug. 23, 2013 at 11:11 PM ET

CHINA
China gets set to launch its first moon lander by year's end

SPACE STATION
Japanese astronaut to command space station in March

The Archon Genomics XPRIZE Presented by Express...

"genome sequencing technology is plummeting in cost and increasing in speed independent of our competition"

"companies can do this for less than $5,000 per genome, in a few days or less — and are moving quickly towards the goals we set for the prize."

# What makes them different?



Much human variation is due to difference in ~ 6 million base pairs (0.1 % of genome) referred to as SNPs

# Single Nucleotide Polymorphism (SNP)

Genomic DNA:

SNP

A
TACATAGCCATCGGTANGTACTCAATGATGATA
G

# From reads to evidence

# From reads to evidence



## I. Comparative

Sequence-wise, individuals of a species are nearly identical

Well curated, annotated "reference" genomes exist



*D. melanogaster, Science,* 2000

*H. sapiens, Nature,* 2000
and *Science,* 2000

*M. musculus, Nature,* 2002



Idea: "Map" reads to their point of origin with respect to a reference, then study differences

# From reads to evidence



## 2. *de novo*

Assume nothing! - let reads tell us everything

Reads with overlapping sequence probably originate from overlapping portions of the subject genome

Encode overlap relationships as a graph



Source: De Novo Assembly Using Illumina Reads. Illumina. 2010



The full genome sequence is a "tour" of the graph

Source: De Novo Assembly Using Illumina Reads. Illumina. 2010
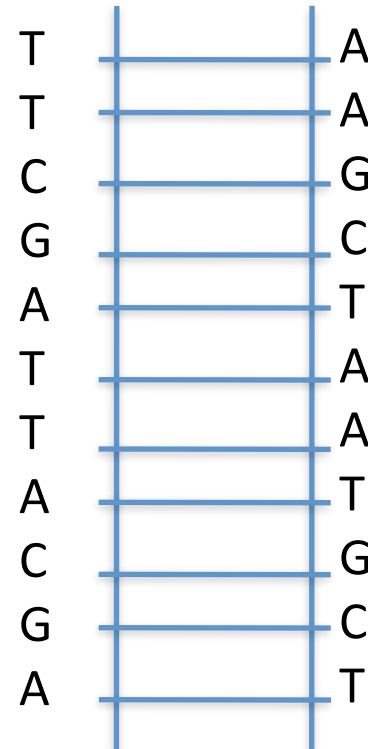http://www.illumina.com/Documents/products/technotes/technote_denovo_assembly.pdf
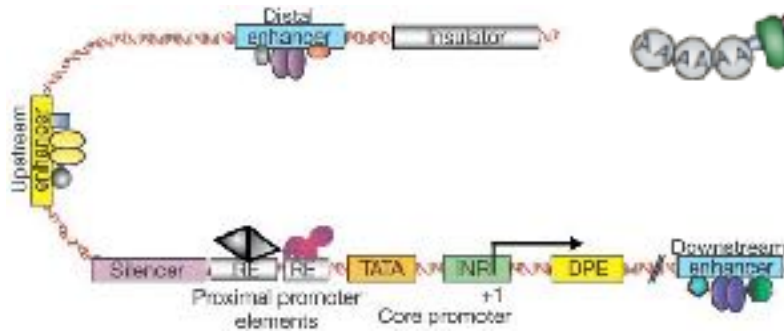
# How many basepair differences?
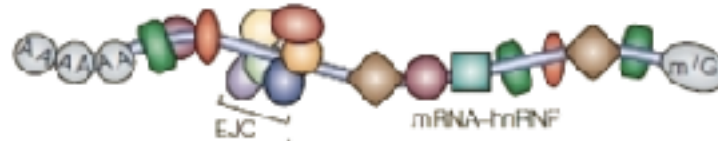
# Computational Biology

Genes encode proteins which are transcribed into mRNA and translated into proteins.
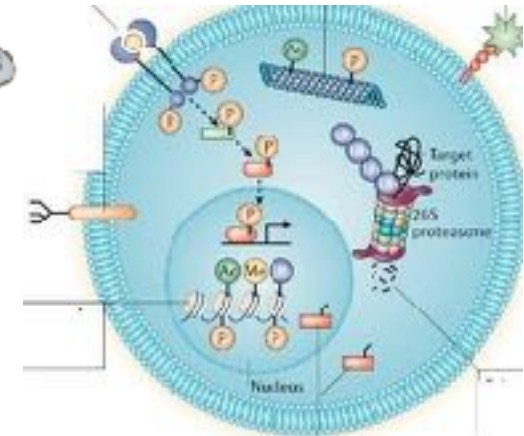
genomics                    transcriptomics                    proteomics



**Major** technological advances allow **unprecedented** data acquisition

# Measurements



Samples (individuals)

DATA MATRIX

# *MLL* translocations specify a distinct gene expression profile that distinguishes a unique leukemia

Scott A. Armstrong[1–4], Jane E. Staunton[5], Lewis B. Silverman[1,3,4], Rob Pieters[6], Monique L. den Boer[6], Mark D. Minden[7], Stephen E. Sallan[1,3,4], Eric S. Lander[5], Todd R. Golub[1,3,4,5]* & Stanley J. Korsmeyer[2,4,8]*

*These authors contributed equally to this work.

# Personal Genomics

# Personal Genomics



Every Genome Tells a Story.

- We need to produce reliable genome measurements, but on much bigger scale (Algorithmics, Systems)

- Multiple genome features, decide which are relevant and significant (Information Retrieval, Data Management)

- Population-based science, interpreted individually (Machine Learning/Statistics, Privacy)

# NHGRI strategic plan

- What does the NIH think genomics should be for the next 10 years?

## PERSPECTIVE

## Charting a course for genomic medicine from base pairs to bedside

Eric D. Green[1], Mark S. Guyer[1] & National Human Genome Research Institute*

There has been much progress in genomics in the ten years since a draft sequence of the human genome was published. Opportunities for understanding health and disease are now unprecedented, as advances in genomics are harnessed to obtain robust foundational knowledge about the structure and function of the human genome and about the genetic contributions to human health and disease. Here we articulate a 2011 vision for the future of genomics research and describe the path towards an era of genomic medicine.

[Nature, Feb. 2011]

# Where do we fit in?

- The major bottleneck in genome sequencing is no longer data generation—the computational challenges around data analysis, display and integration are now rate limiting. New approaches and methods are required to meet these challenges.

- **Data analysis**
  - Computational tools are quickly becoming inadequate for analysing the amount of genomic data that can now be generated, and this mismatch will worsen. Innovative approaches to analysis, involving close coupling with data production, are essential.

- **Data integration**
  - Genomics projects increasingly produce disparate data types (for example, molecular, phenotypic, environmental and clinical), so computational approaches must not only keep pace with the volume of genomic data, but also their complexity. New integrative methods for analysis and for building predictive models are needed.

- **Visualization**
  - In the past, visualizing genomic data involved indexing to the one-dimensional representation of a genome. New visualization tools will need to accommodate the multidimensional data from studies of molecular phenotypes in different cells and tissues, physiological states and developmental time. Such tools must also incorporate non-molecular data, such as phenotypes and environmental exposures. The new tools will need to accommodate the scale of the data to deliver information rapidly and efficiently.

- **Computational tools and infrastructure**
  - Generally applicable tools are needed in the form of robust, well-engineered software that meets the distinct needs of genomic and non-genomic scientists. Adequate computational infrastructure is also needed, including sufficient storage and processing capacity to accommodate and analyse large, complex data sets (including metadata) deposited in stable and accessible repositories, and to provide consolidated views of many data types, all within a framework that addresses privacy concerns. Ideally, multiple solutions should be developed[105].

# Where do we fit in?

- Meeting the computational challenges for genomics requires scientists with expertise in biology as well as in informatics, computer science, mathematics, statistics and/or engineering.

- *A new generation of investigators who are proficient in two or more of these fields must be trained and supported.*

# What else is the class about?

- Gives you an example of end-to-end use of what you've learned as CS as a practice
  - We discuss the design and analysis of algorithms (e.g., string algorithms, dynamic programming, iterative optimization methods)
  - We implement algorithms (python)
  - We analyze data (also in python)
- We also learn about biology, medicine and why government shutdowns are really awful

# Administrative Details

Class webpage:

1. http://www.hcbravo.org/cmsc423

Everything you want to know is there.

Todo after class today

1) Enroll in Piazza class

2) Enroll in Rosalind pre-lecture and final submission pages (links posted in Piazza)

3) Complete course survey on ELMS

For next class

1) Reading

2) Pre-lecture reading quiz on ELMS

3) Pre-lecture rosalind exercises