

# INTRO TO INTRO

---

# What is Computational Biology?

Using

computation  
algorithms  
mathematics  
statistics  
visualization  
databases  
simulations

to help understand

biology  
molecules  
reactions  
cells  
systems  
disease  
evolution

In Bio131, we will use **algorithms** to  
help understand **biological sequences**

# The Quest to Sequence the Human Genome



# The Quest to Sequence the Human Genome

*“Regardless of who wins this sprint, the next race--to make sense of the genome--will be a marathon with many runners...Just the first step, which researchers call annotation, could very well take many months.” – Karow, 2000*

## ENCyclopedia Of DNA Elements (ENCODE) Project

2003 ENCODE Project Established

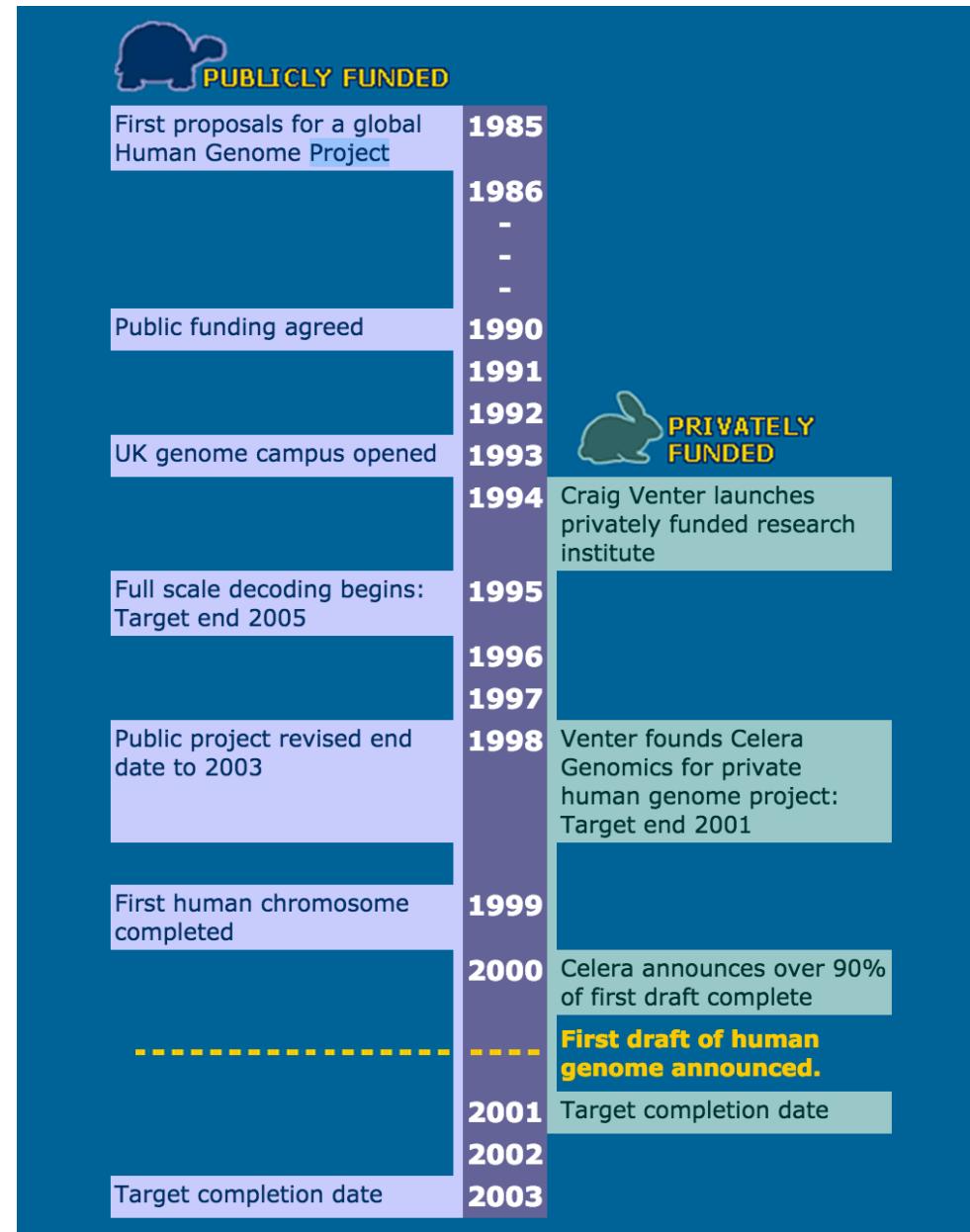
2007 ENCODE Pilot Completed

2012 ENCODE Project Completed

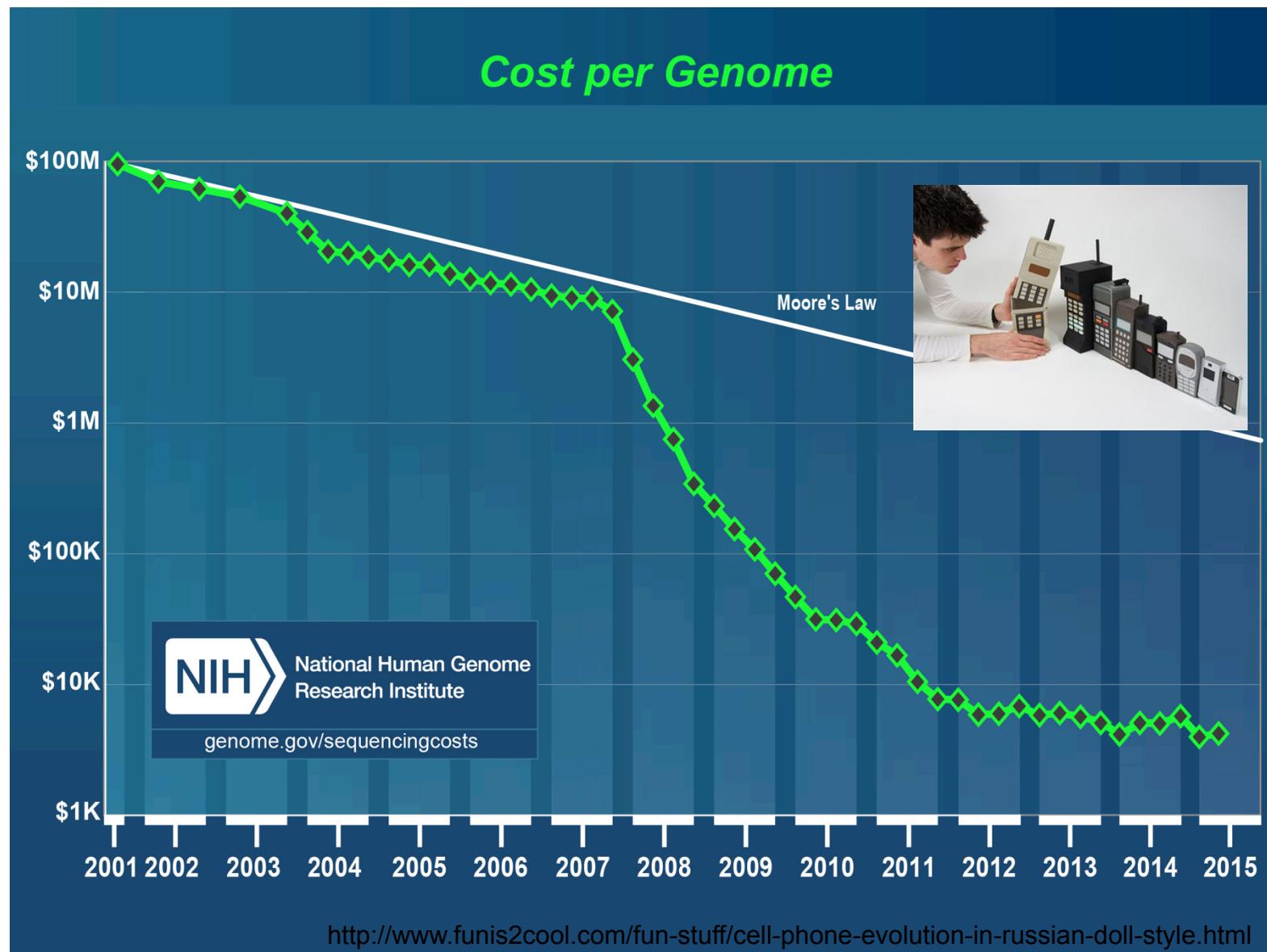
[The Human Genome Race](#), Julia Karow, Scientific American, 2000.

<http://www.genome.gov/10001763>

[http://news.bbc.co.uk/hi/english/static/in\\_depth/sci\\_tech/2000/human\\_genome/timeline/default.stm](http://news.bbc.co.uk/hi/english/static/in_depth/sci_tech/2000/human_genome/timeline/default.stm)



# A Deluge of Data



# A Deluge of Data

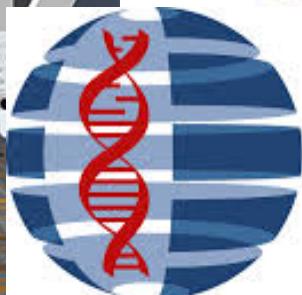
**1000 Genomes**  
A Deep Catalog of Human Genetic Variation



**1000 Plant Genomes**



**Fish-T1K**  
Transcriptomes of 1000 Fishes



**International Cancer Genome Consortium**



**NIH HUMAN MICROBIOME PROJECT**



**THE CANCER GENOME ATLAS**

National Cancer Institute  
National Human Genome Research Institute

# Example: Teleost Fishes



# Example: Teleost Fishes

## SCIENTIFIC DATA

A graphic consisting of a 5x6 grid of binary digits (0s and 1s) in blue and green, arranged in a roughly triangular shape that tapers to the right.**OPEN**

### Data Descriptor: Whole genome sequencing data and *de novo* draft assemblies for 66 teleost species

Martin Malmstrøm<sup>1</sup>, Michael Matschiner<sup>1</sup>, Ole K. Tørresen<sup>1</sup>, Kjetill S. Jakobsen<sup>1</sup>  
& Sissel Jentoft<sup>1,2</sup>

Received: 14 September 2016

Accepted: 07 December 2016

Published: 17 January 2017

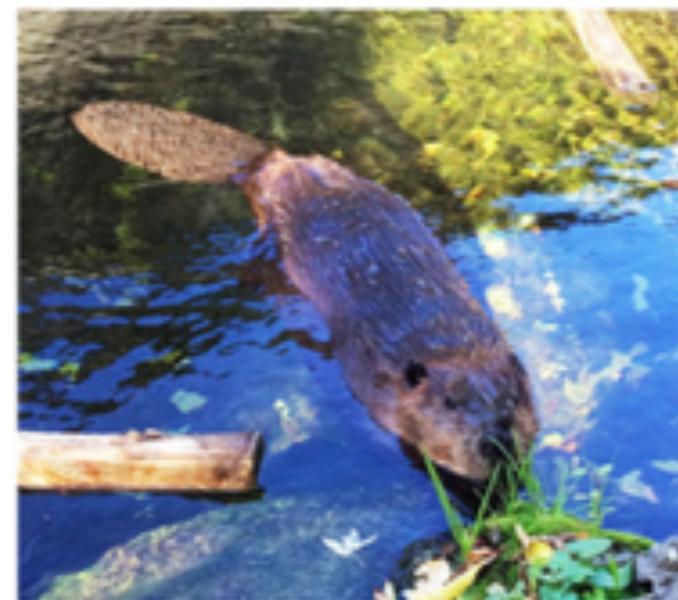
Teleost fishes comprise more than half of all vertebrate species, yet genomic data are only available for 0.2% of their diversity. Here, we present whole genome sequencing data for 66 new species of teleosts, vastly expanding the availability of genomic data for this important vertebrate group. We report on *de novo* assemblies based on low-coverage (9–39  $\times$ ) sequencing and present detailed methodology for all analyses. To facilitate further utilization of this data set, we present statistical analyses of the gene space completeness and verify the expected phylogenetic position of the sequenced genomes in a large mitogenomic context. We further present a nuclear marker set used for phylogenetic inference and evaluate each gene tree in relation to the species tree to test for homogeneity in the phylogenetic signal. Collectively, these analyses illustrate the robustness of this highly diverse data set and enable extensive reuse of the selected phylogenetic markers and the genomic data in general. This data set covers all major teleost lineages and provides unprecedented opportunities for comparative studies of teleosts.

# Example: The Canadian Beaver

## De Novo Genome and Transcriptome Assembly of the Canadian Beaver (*Castor canadensis*)

Si Lok,<sup>\*,†,1</sup> Tara A. Paton,<sup>\*,†</sup> Zhuozhi Wang,<sup>\*,†</sup> Gaganjot Kaur,  
Wilson W. L. Sung,<sup>\*,†</sup> Joseph Whitney,<sup>\*,†</sup> Janet A. Buchanan,<sup>\*</sup>  
Beverly Apresto,<sup>\*,†</sup> Nan Chen,<sup>\*,†</sup> Matthew Coole,<sup>\*,†</sup> Travis J. |  
Sanjeev Pullenayegum,<sup>\*,†</sup> Kozue Samler,<sup>\*,†</sup> Arun Shipstone,<sup>\*</sup>  
Sergio L. Pereira,<sup>\*,†</sup> Pirooz Rostami,<sup>\*,†</sup> Carol Ann Ryan,<sup>\*,†</sup> Ar  
Yogi Sundaravadanam,<sup>§</sup> Jared T. Simpson,<sup>§,\*\*</sup> Burton K. Lim,  
Christopher J. Dutton,<sup>††</sup> Kevin C. R. Kerr,<sup>††</sup> Maria Franke,<sup>††</sup>  
and Stephen W. Scherer<sup>§§,\*\*\*\*,†</sup>

\*The Centre for Applied Genomics and <sup>†</sup>Program in Genetics and Genomics, University of Toronto, Toronto, Ontario M5G 0A4, Canada, <sup>§§</sup>McLaughlin Centre, Faculty of Medicine, University of Waterloo, Waterloo, Ontario N2L 3G1, Canada, <sup>††</sup>Donnelly Centre for Cellular and Biomolecular Research, University of Toronto, Toronto, Ontario M5S 1A8, Canada, <sup>\*\*</sup>Department of Medical Biophysics, University of Toronto, Toronto, Ontario M5S 1A8, Canada, <sup>††</sup>Royal Ontario Museum, Toronto, Ontario M5S 2C6, Canada, <sup>†††</sup>Toronto Zoo, Toronto, Ontario M1B 5K7, Canada, and <sup>\*\*\*\*</sup>Department of Molecular Genetics, Faculty of Medicine, University of Toronto, Toronto, Ontario M5S 1A8, Canada



to,

for

to,

# Example: An Ancient Genome [Spring 2016]



**Neolithic woman from Ballynahatty, Ireland (~5,000 years ago)**

- DNA suggests Near Eastern origin

**Bronze-age individuals from similar region (~4,000 years ago)**

- DNA suggests European origin

# Example: The Octopus Genome [Fall 2015]



## Compared to Human:

⬇ Genome Size 2.7Gb vs. 3.0Gb

⬆ # of Protein Coding Genes 33K vs. 25K

⬆ # of Genes Related to Neuron Regulation 168 vs. 70

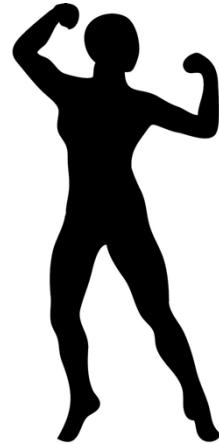
⬇ # of Neurons 500M vs. ~20B (mouse has 4M)

# Overview of Topics

- **Topic 1:** How do cells “read” the genome?
- **Topic 2:** Where do DNA-binding proteins bind?
- **Topic 3:** How do we compare biological sequences?
- **Topic 4:** How do we assemble genomes?
- **Topic 5:** How do we compare entire genomes?

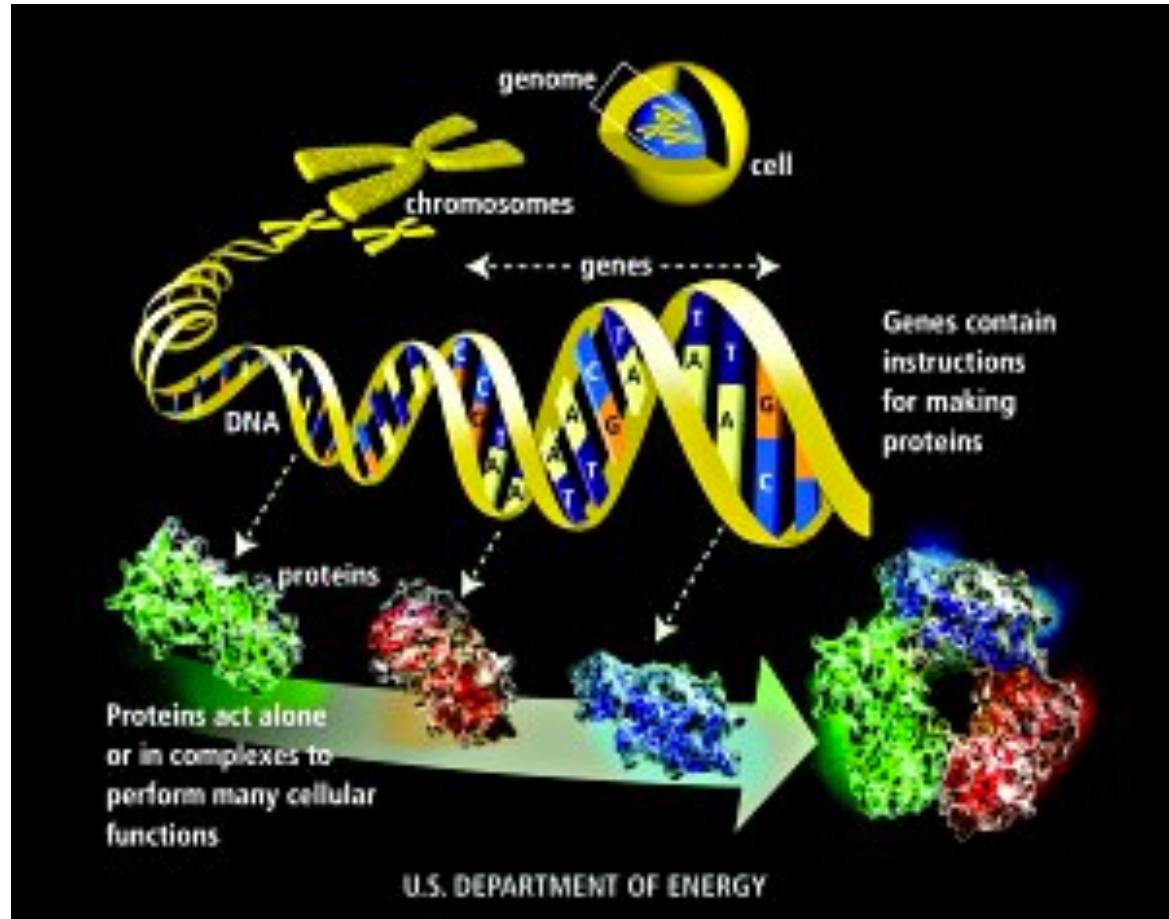
*All of these topics involve  
biological sequence analysis!*

# Topic 1: How do Cells “Read” the Genome?



$$3.72 \times 10^{13} \text{ cells} = \\ 37,200,000,000,000 \\ = 37.2 \text{ trillion cells}$$

Each **cell** has the same genome with **3 billion bases** = 3,000,000,000



# Topic 1: Program the Central Dogma

## Input

AGTGCTCGACGCAGCTACGACTACGTTATCGA...

## Intermediate Steps

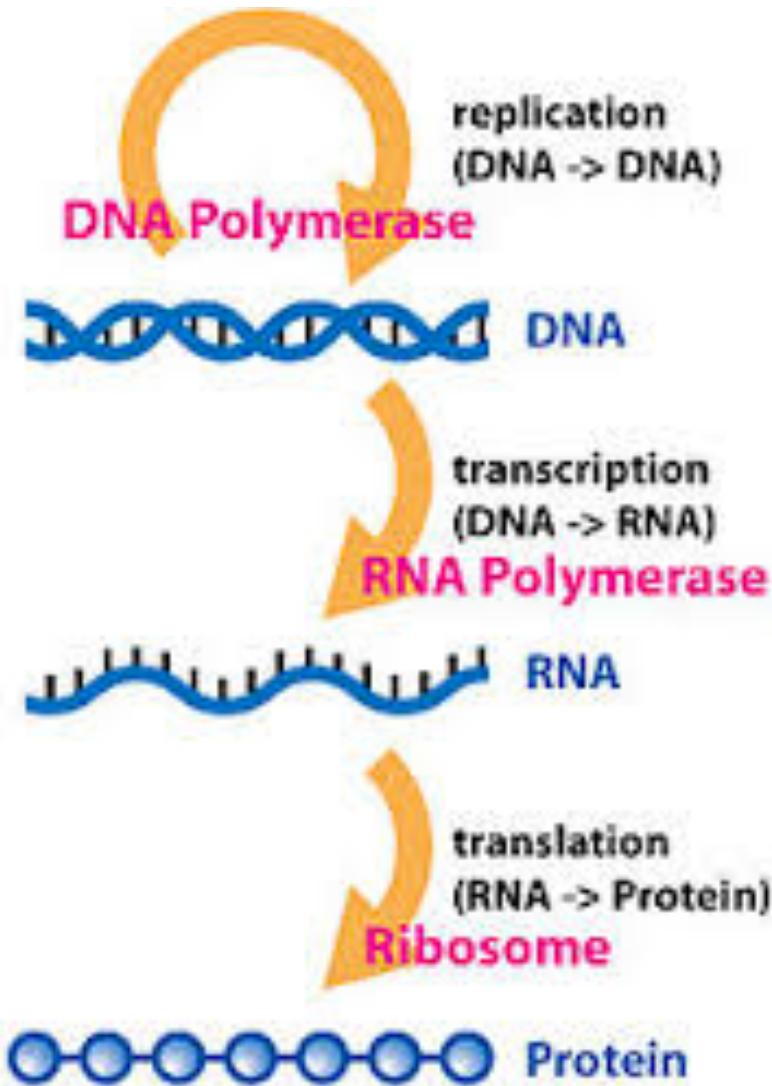
AGUGCUCGACGCAGCUACGACUACGUUAUCGA...

CUCGACGCA CUACGU GA...

CUCGACGCACUACGUGA...

## Output

KTYSILP...

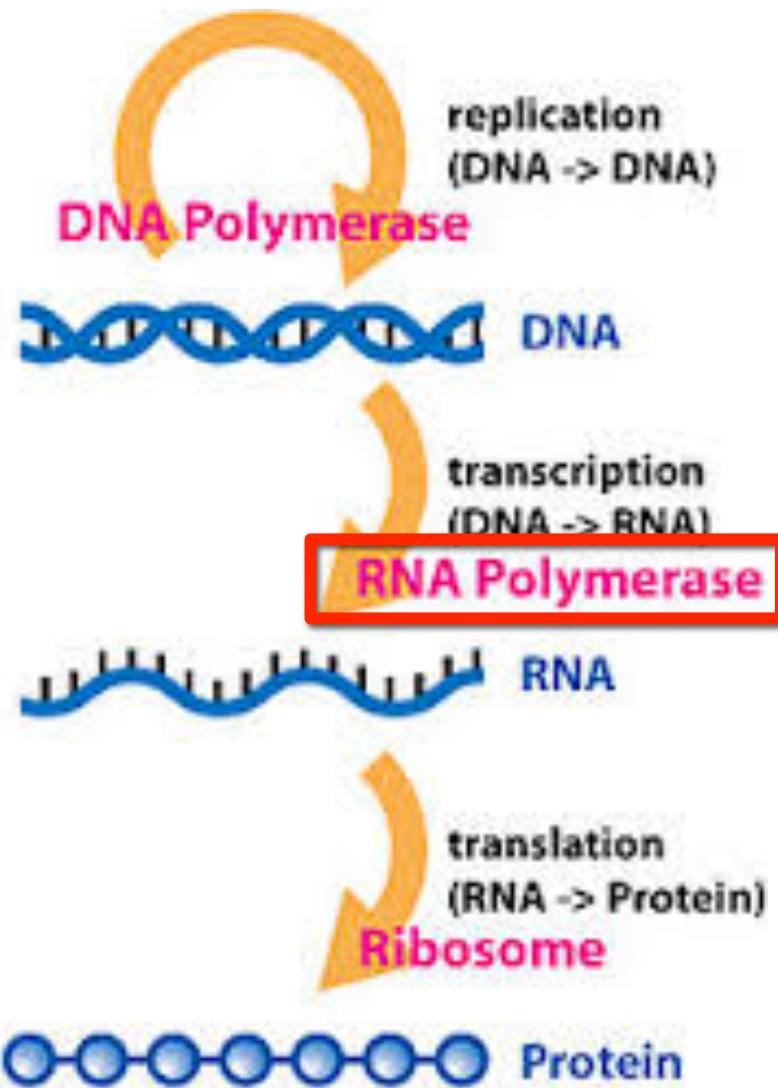
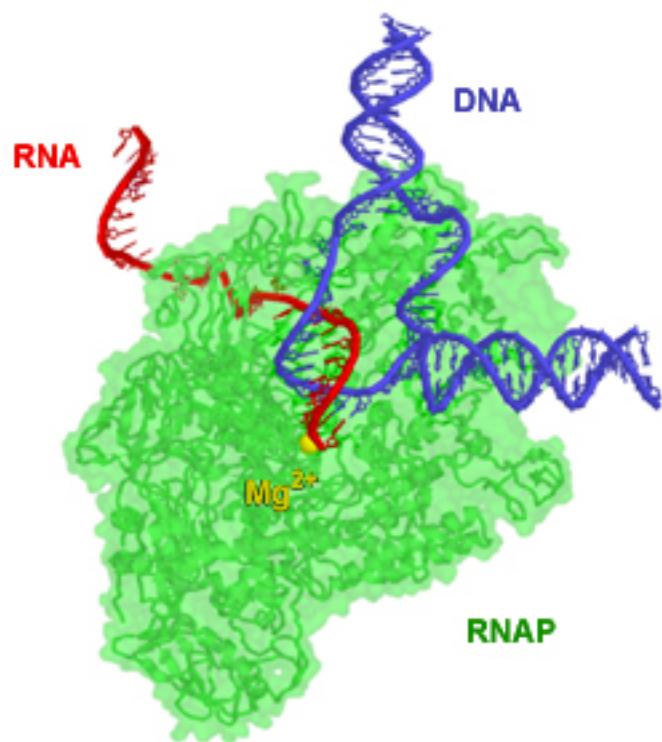


# By the end of...

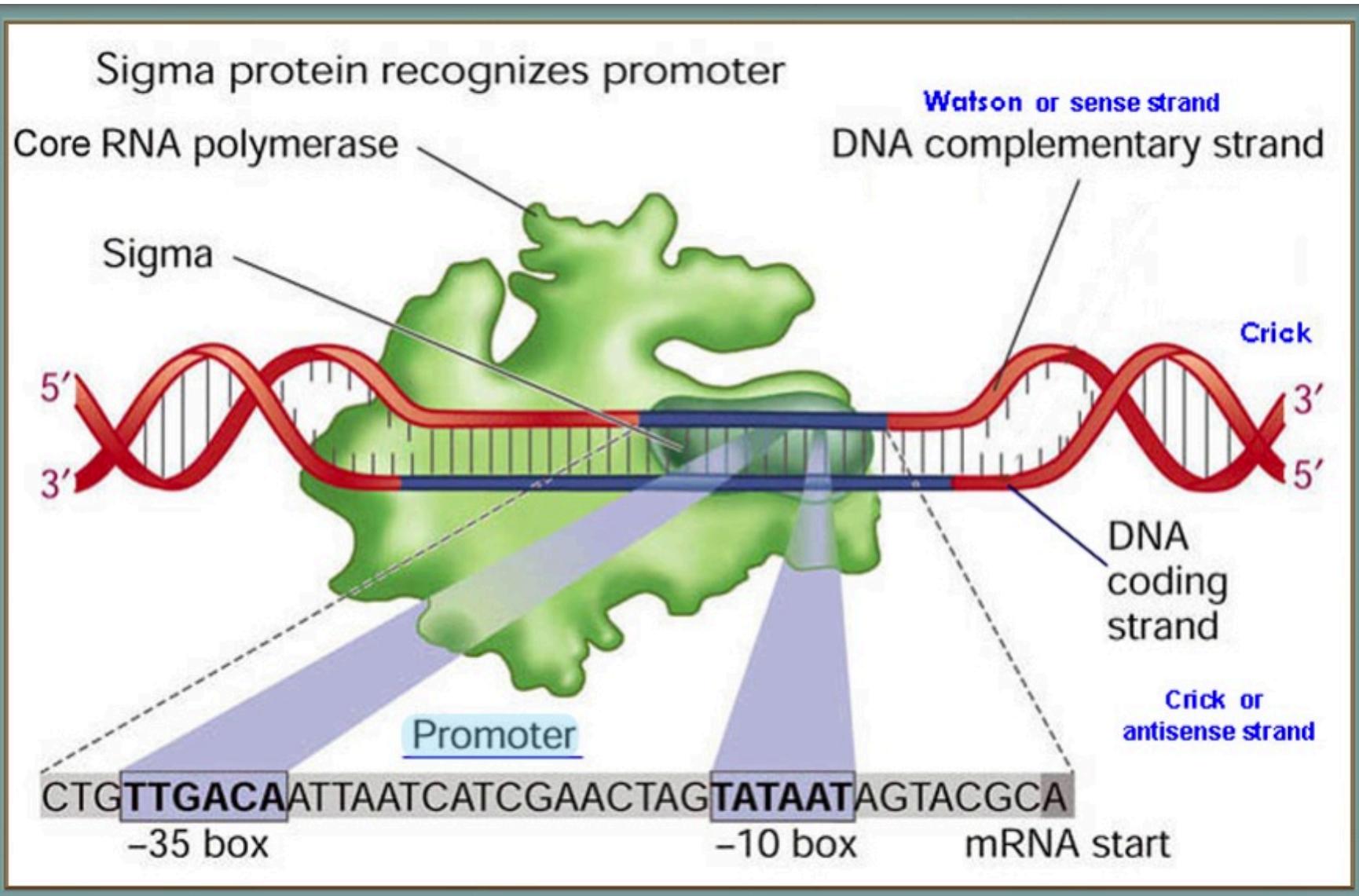
Topic 1:

Write a Python program to convert **any** gene to its protein sequence.

# Topic 2: Where do DNA-binding Proteins Bind?



# Topic 2: Where do DNA-binding Proteins Bind?



# By the end of...

Topic 1:

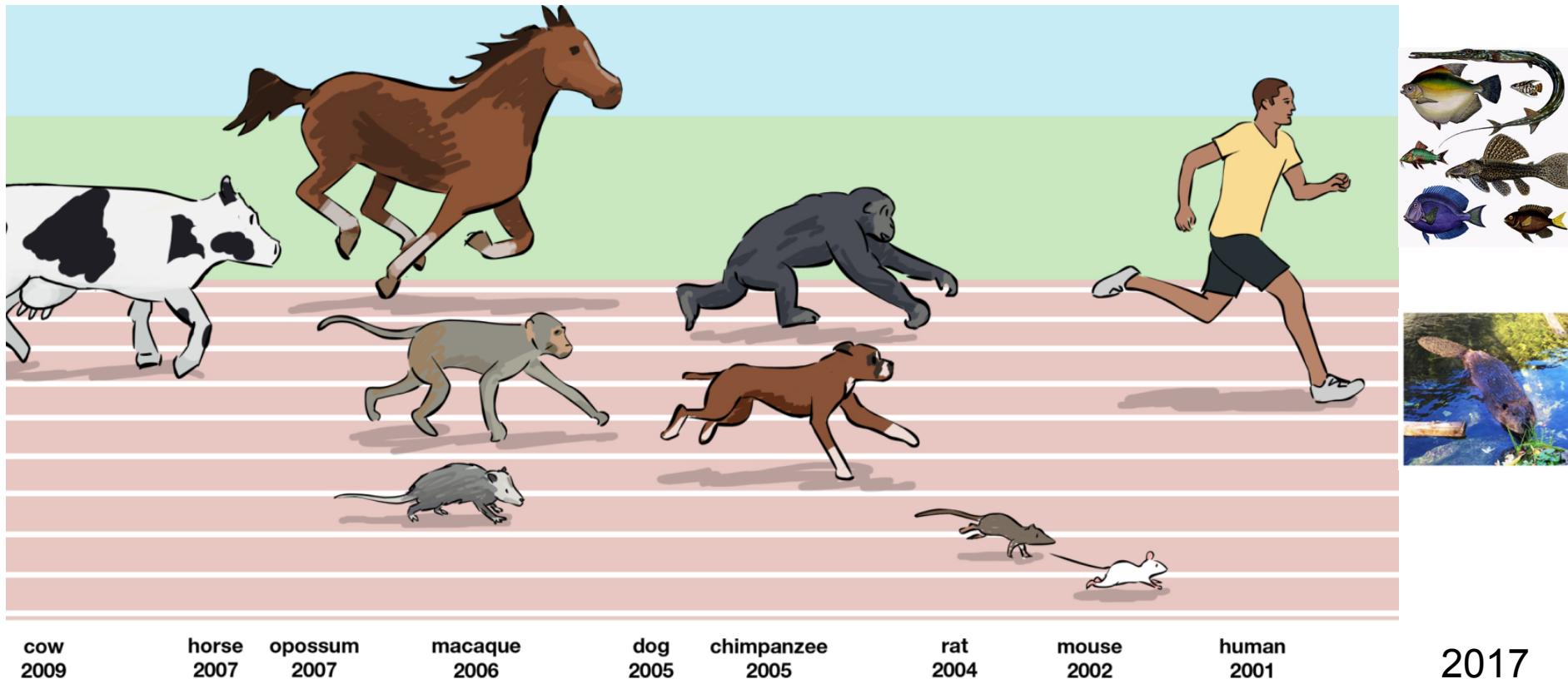
Write a Python program to convert **any** gene to its protein sequence.

Topic 2:

Write a Python program to identify and summarize motifs from a collection of sequences.

# Topic 5: Assembling Genomes “From Scratch”

- How was the first human genome assembled?
- How are the genomes of other species assembled?

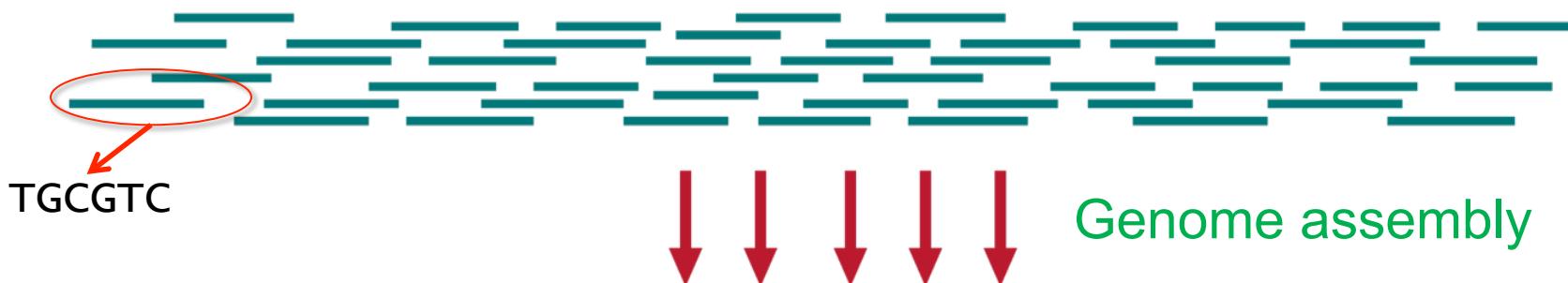


# Topic 5: Assembling Genomes “From Scratch”

Multiple (unsequenced) genome copies



Sequenced Reads



Assembled genome

...GGCATGCGTCAGAAACTATCATAGCTAGATCGTACGTAGCC...

# By the end of...

Topic 1:

Write a Python program to convert **any** gene to its protein sequence.

Topic 2:

Write a Python program to identify and summarize motifs from a collection of sequences.

Topic 3:

Write a Python program to simulate an assembler

# Topic 4: How do we compare sequences?

(hypothetical)

Human Gene: ACTCGACTGAGAGGATTCGAGCATGA

Mouse Gene: ACTCAACTGAGA**TTCGAGCTTCA**ATGA

ACTCGACTGAGA <b>GG</b> ATTC <b>GAGC</b> ATGA	
ACTCAACTGAGA <b>TT</b> CGAGC <b>TTCA</b> ATGA	

11 substitutions

If we “shift” the sequences by adding gaps,  
can we improve the number of substitutions?

Finding an Optimal  
Sequence Alignment

# By the end of...

Topic 1:

Write a Python program to convert **any** gene to its protein sequence.

Topic 2:

Write a Python program to identify and summarize motifs from a collection of sequences.

Topic 3:

Write a Python program to simulate an assembler

Topic 4:

Write multiple Python programs to align pairs of DNA or protein sequences.

# Topic 4: How do we compare genomes?

## Human Genome Project

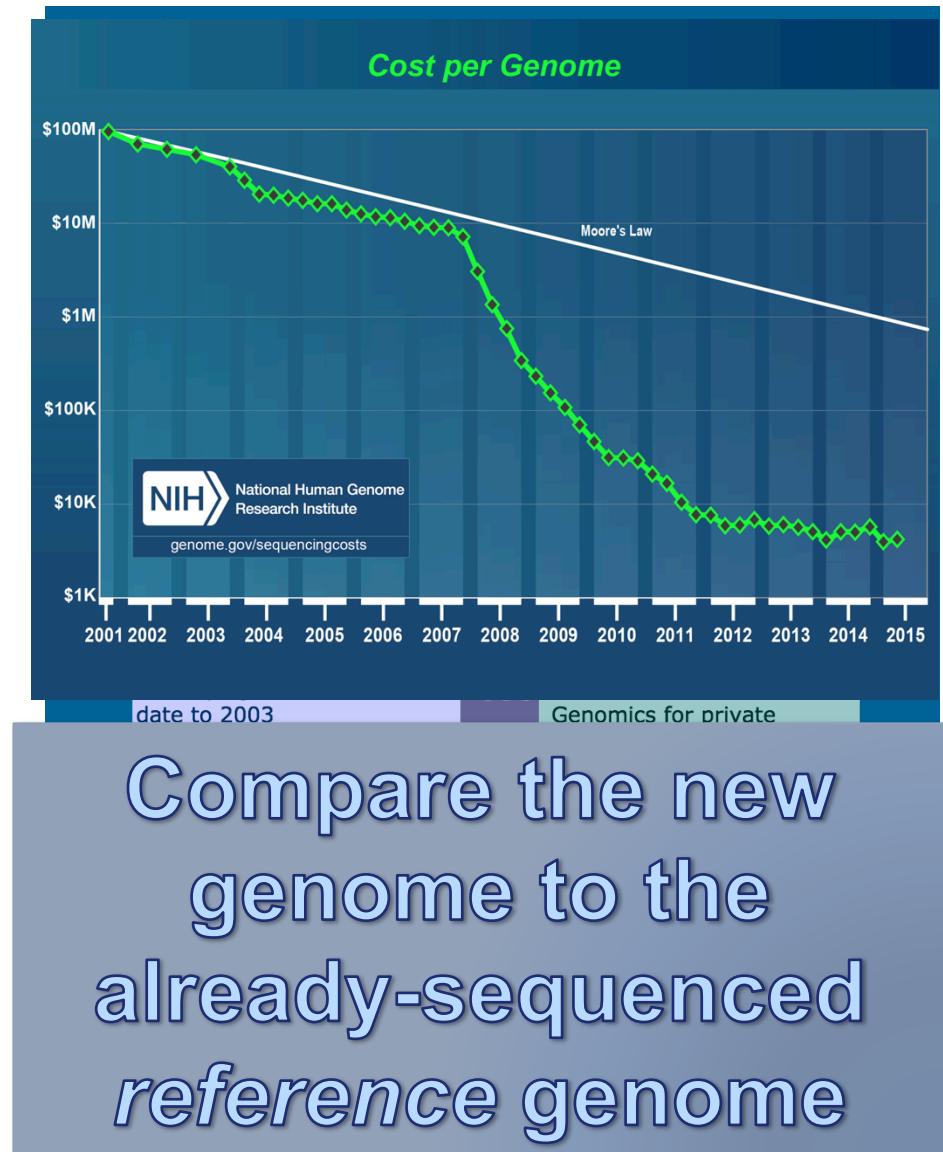
- \$2.7 billion
- About 6 years

## Celera Genomics

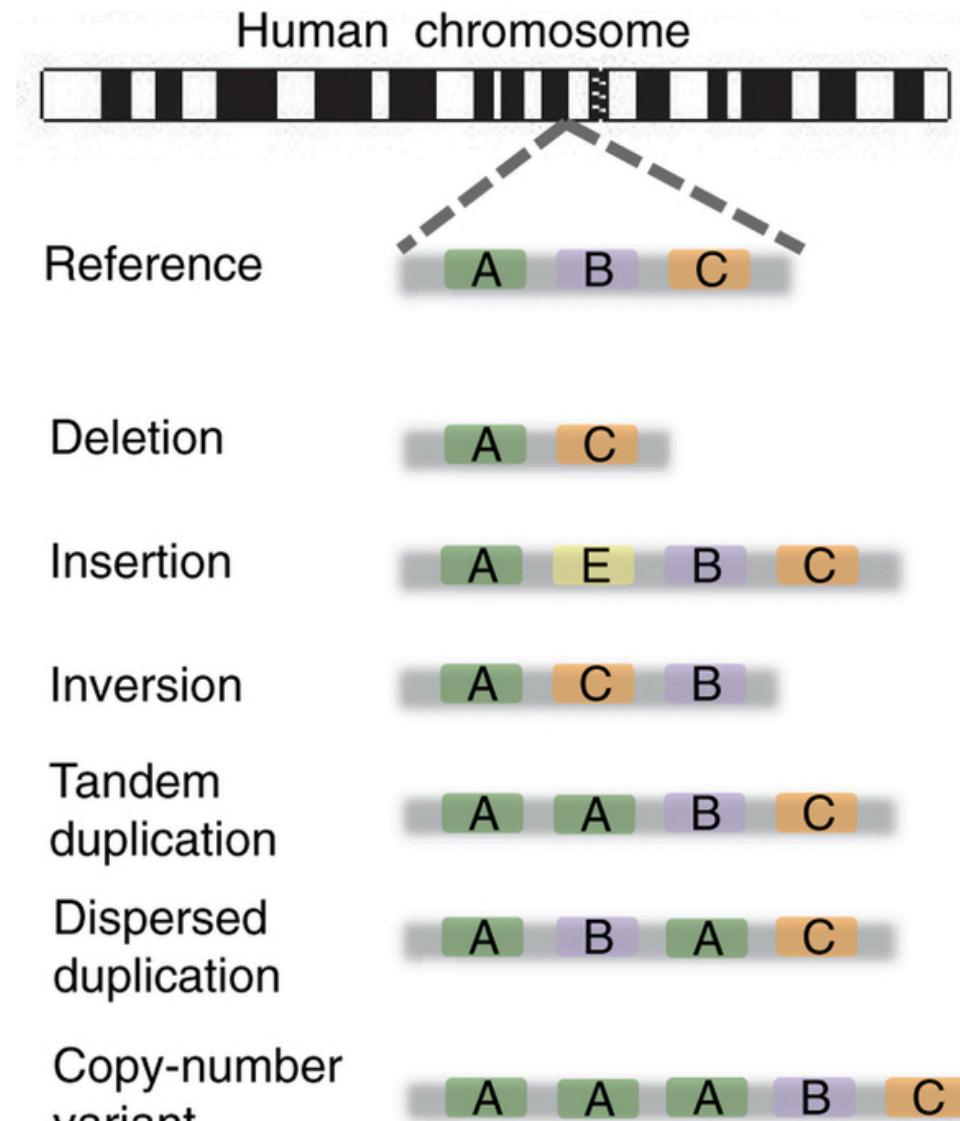
- \$300 million
- About 3 years

## Today:

- \$1,000 (announced 2015)
- \$4,000 typically
- A week (I'm generous)



# Topic 5: Genome Rearrangements



# By the end of...

Topic 1:

Write a Python program to convert **any** gene to its protein sequence.

Topic 2:

Write a Python program to identify and summarize motifs from a collection of sequences.

Topic 3:

Write a Python program to simulate an assembler

Topic 4:

Write multiple Python programs to align pairs of DNA or protein sequences.

Topic 5:

Visualize large genomic rearrangements between two genomes.

# By the end of...

Topic 1:

Write a Python program to convert **any** gene to its protein sequence.

Topic 2:

Write a Python program to identify and summarize motifs from a collection of sequences.

Topic 3:

Write a Python program to simulate an assembler

Topic 4:

Write multiple Python programs to align pairs of DNA or protein sequences.

Topic 5:

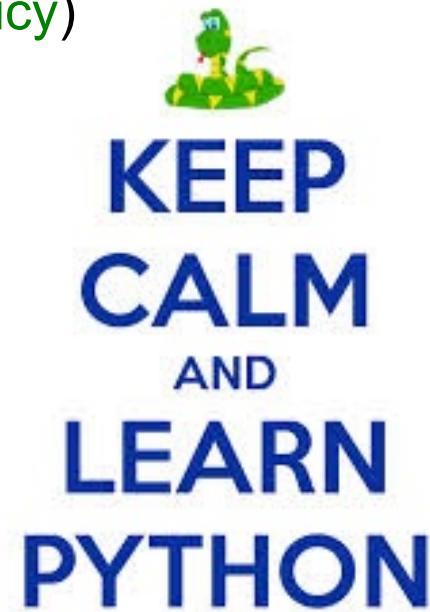
Visualize large genomic rearrangements between two genomes.

Topic 6:

Student-Selected Topics & Final Projects

# About This Class

- Hybrid biology/computer science class
- We will learn Python!
- Course Components
  - 8 graded programming assignments ([Collaboration Policy](#))
  - 2 in-class exams on biology/computational concepts
  - 2 code conferences
  - Labs handed in for participation points
  - Final project on a topic of your choice
- More information on the syllabus



# About This Class

- Moodle & Website

- <https://moodle.reed.edu/course/view.php?id=1251>
- <http://reed.edu/biology/courses/bio131/>

- Textbook

- *Bioinformatics Algorithms: an Active Learning Approach*
- Edition 2, Volume 1
- Pevzner & Compeau
- [www.bioinformaticsalgorithms.org/](http://www.bioinformaticsalgorithms.org/)



- Additional problems will be assigned via **ROSALIND**

- <http://www.rosalind.info/>

# About This Class

- Respect different levels of experience
- Adhere to the **Collaboration Policy**
  - Clearly defines the Honor Principle as it relates to coursework
  - Properly “cite” people you worked with
- **Everyone** will have an oh sh!t moment
- On the flip side, real life is more important than class
  - If something is happening, **take time** to deal with it
  - Come to me with any issues, but be aware that I am an **obligated reporter**
- **Save, Save, Save, Save, Save**



# People

## Professor

Anna Ritz

[aritz@reed.edu](mailto:aritz@reed.edu), B200B

*Office Hours*

*W/Th 11am-12pm  
or by appointment*



Anna



Mina (Math)

## Tutors/TAs

Rose Driscoll (Mon Lab; Thurs DoJo)

Mina Marden (Tues Lab)

Elaine Kushkowski (Tues DoJo)



Elaine (ES-Bio)



Rose (Alt Bio)

# HW1 Out (Moodle)



- Excel assignment
  - Excel is available for free on the CIS software page
- Due Monday Jan 30 before class
- Part 1: Warmup with Excel Functions
- Part 2: Sequence Analysis (Wed. lecture)

Enter numbers/strings in green cells

Enter formulas (starting with "=") in blue cells.

Enter responses/text in red cells

# Upcoming Schedule

<b>Week 1</b>	
Mon 1/23	Lecture HW1 Out
<b>Lab1: Python Setup &amp; Command Line Tools</b>	
Wed 1/25	Lecture
Fri 1/27	Lecture
<b>Week 2</b>	
Mon 1/30	Lecture HW1 Due HW2 Out
<b>Lab2: Python Pattern Maker</b>	
Wed 2/1	Lecture
Fri 2/3	Lecture