



CS313 Computational Biology



Course Information

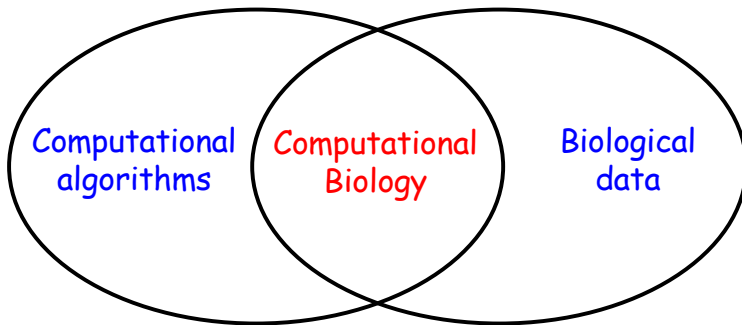
Instructor: Brian Tjaden

Pronouns: He, Him, His

Course Materials: <http://cs.wellesley.edu/~cs313>



Computational Biology is Multidisciplinary



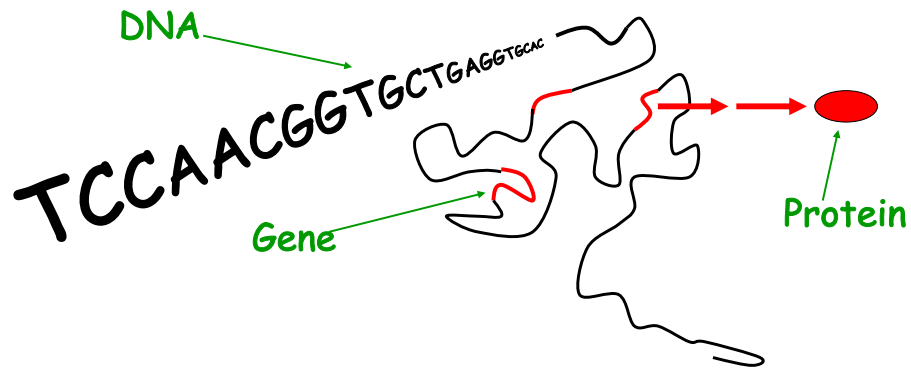
DNA the Molecule of Life

The infographic illustrates the hierarchy of genetic information. At the top right is a yellow sphere representing a **cell**. Below it are yellow X-shaped structures representing **chromosomes**. In the center is a DNA double helix structure with the label **DNA**. A specific segment of the DNA is highlighted with a bracket and labeled **gene**. The bases of the gene are shown as colored blocks with letters: C (blue), G (orange), A (yellow), and T (purple).

Y:GG 00-0481



DNA: simplified



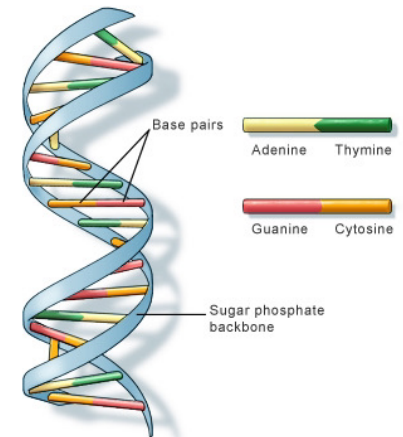
DNA: "program" for cell processes
Proteins: execute cell processes

A-5



DNA Structure

- Double helix
- Deoxyribose (sugar) - phosphate backbone
- Four bases - A, T, G, C
- Base pairing



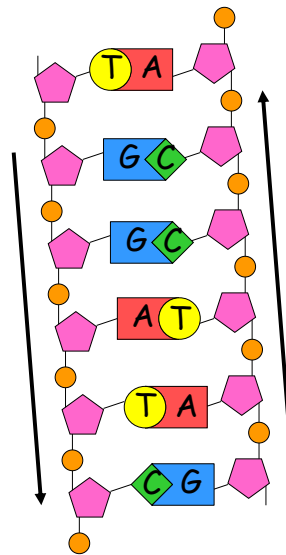
U.S. National Library of Medicine

A-6



DNA Structure

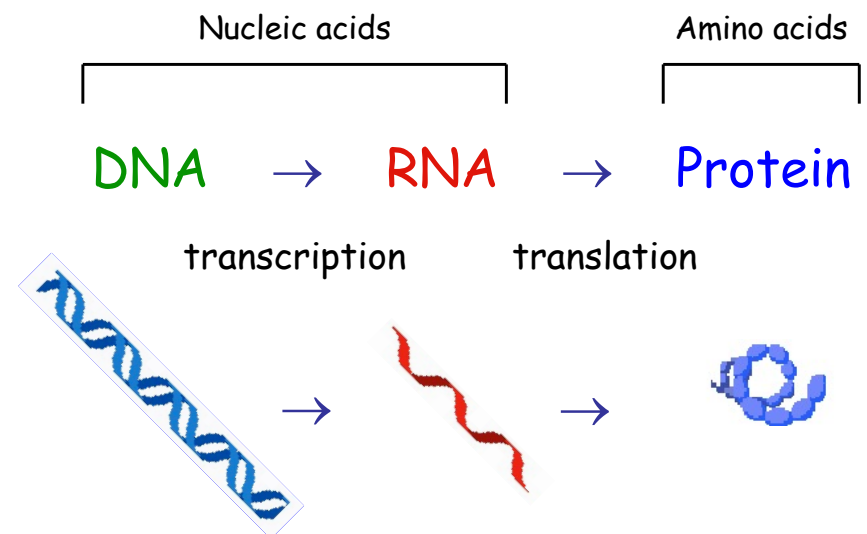
- **Information polarity** (anti-parallel strands)
- **Either strand can function as a template** (complementary strands)



A-7



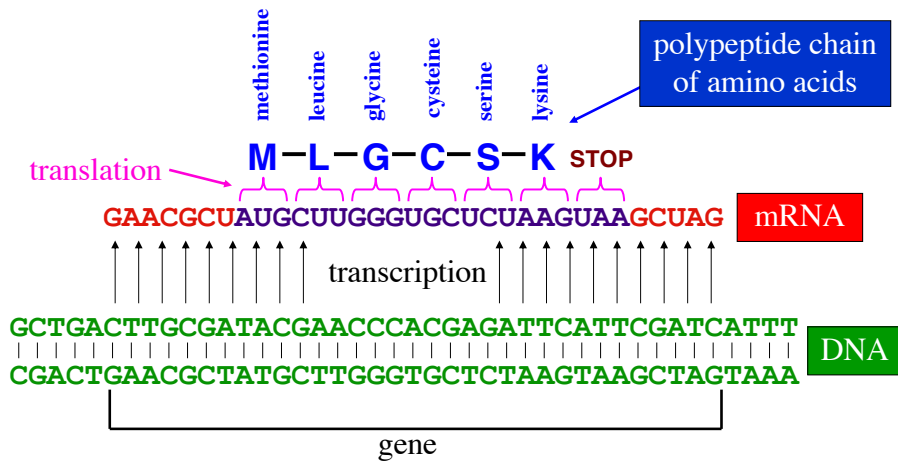
Information Flow



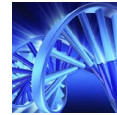
A-8



DNA → RNA → Protein



A - 9



The Genetic Code

- 61 amino acid codons
- 1 start codon (Met)
- 3 stop codons

| | | Second letter | | | | |
|---|-----------|---------------|------------|------------|------------------|--|
| | | U | C | A | G | |
| U | UUU } Phe | UCU } Ser | UAU } Tyr | UGU } Cys | U C A G | |
| | UUC } Leu | UCC } Ser | UAC } Stop | UGC } Cys | | |
| | UUA } Leu | UCA } Stop | UAA } Stop | UGA } Stop | | |
| | UUG } Leu | UCG } Ser | UAG } Stop | UGG } Trp | | |
| C | CUU } Leu | CCU } Pro | CAU } His | CGU } Arg | U C A G | |
| | CUC } Leu | CCC } Pro | CAC } His | CGC } Arg | | |
| | CUA } Leu | CCA } Pro | CAA } Gln | CGA } Arg | | |
| | CUG } Leu | CCG } Pro | CAG } Gln | CGG } Arg | | |
| A | AUU } Ile | ACU } Thr | AAU } Asn | AGU } Ser | U C A G | |
| | AUC } Ile | ACC } Thr | AAC } Asn | AGC } Ser | | |
| | AUA } Met | ACA } Thr | AAA } Lys | AGA } Arg | | |
| | AUG } Met | ACG } Thr | AAG } Lys | AGG } Arg | | |
| G | GUU } Val | GCU } Ala | GAU } Asp | GGU } Gly | U C A G | |
| | GUC } Val | GCC } Ala | GAC } Asp | GGC } Gly | | |
| | GUA } Val | GCA } Ala | GAA } Glu | GGA } Gly | | |
| | GUG } Val | GCG } Ala | GAG } Glu | GGG } Gly | | |

A - 10



The Genetic Code

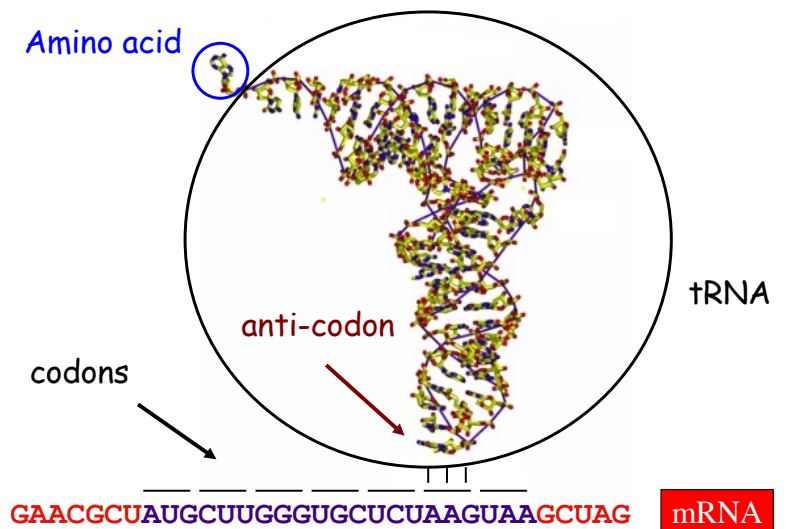
- Triplet code
- Non-overlapping codons
- Start and stop codons
- Degeneracy

4 nucleotides, 20 amino acids

A - 11



tRNAs (transfer RNAs)

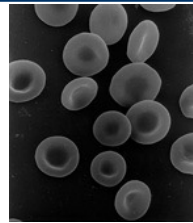


A - 12

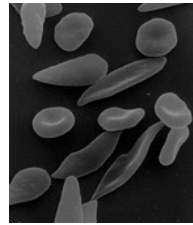


Mutations

- Changes in DNA occur, despite cell's best efforts
- Spontaneous events, copying errors, enviromental factors
- Mutations might change gene function
- Can be harmful, neutral, or beneficial



Normal RBCs



Sickle cell anemia

A - 13

```

ATGGTGCAIC TGACTCCTGA GGAGAAGTCT GCCGTTACTG CCCTGTGGGG CAAGGTGAAC GTGGATGAAG TTGGTGGTGA
GGCCCTGGGC AGGCTGCTGG TGGTCTACCC TTGGACCCAG AGGTTCTTTG AGTCCTTTGG GGATCTGTCC ACTCCTGATG
CTGTTAATGG CAACCTAAG GTGAAGGCTC ATGGCAAGAA AGTGTCTGGT GCCTTTAGTG ATGGCCTGGC TCACCTGGAC
AACCTCAAGG GCACCTTTGC CACACTGAGT GAGCTGCACT GTGACAAGCT GCACGTGGAT CCGAGACT TCAGGCTCCT
GGGCAACGTG CTGGTCTGTG TGCTGGCCCA TCACTTTGGC AAAGAATTCA CCCCACCAGT GCAGGCTGCC TATCAGAAAG
TGTTGGCTGG TGTGGCTAAT GCCTGGCCC ACAAGTATCA CTA

```



Ultraviolet (UV) Light Causes Sunburn and DNA Damage



A - 14

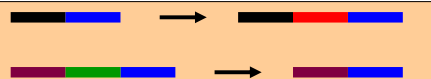


Types of Mutations

Single base substitutions

A → T

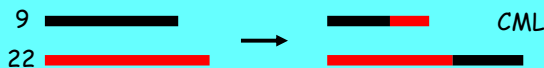
Insertions and Deletions



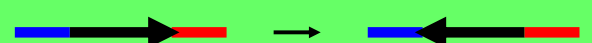
Amplifications



Translocations



Inversions



Sample Genomes

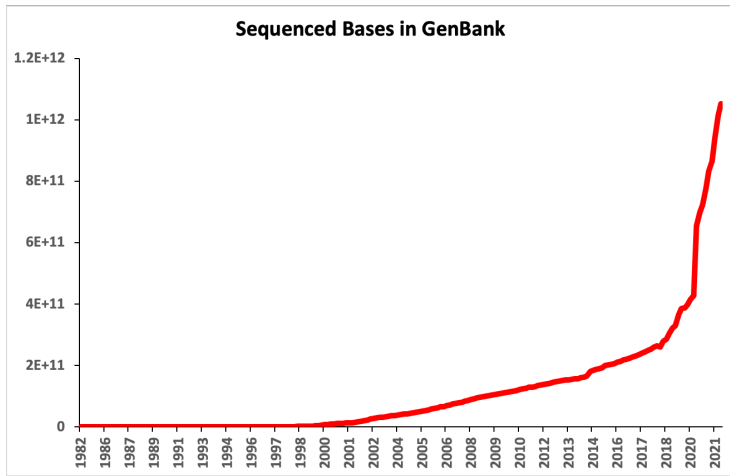
A genome is the total DNA in a cell

| Species | Genome Size | # of Genes |
|--------------------------------|--------------|------------|
| SARS-CoV-2 | 30 thousand | |
| Epstein-Barr virus | 172 thousand | |
| <i>Escherichia coli</i> | 4.6 million | |
| <i>Drosophila melanogaster</i> | 122 million | |
| <i>Homo sapiens</i> | 3.3 billion | |
| <i>Paris japonica</i> | 150 billion | |

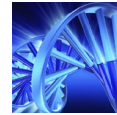
A - 16



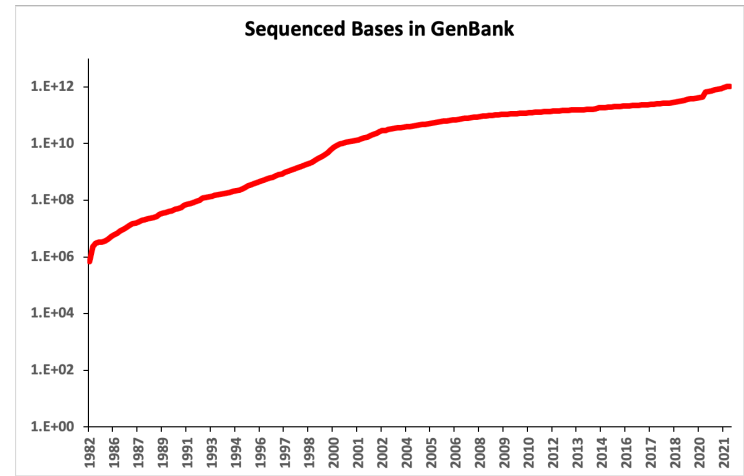
GenBank Growth



A - 17



GenBank Growth (log scale)

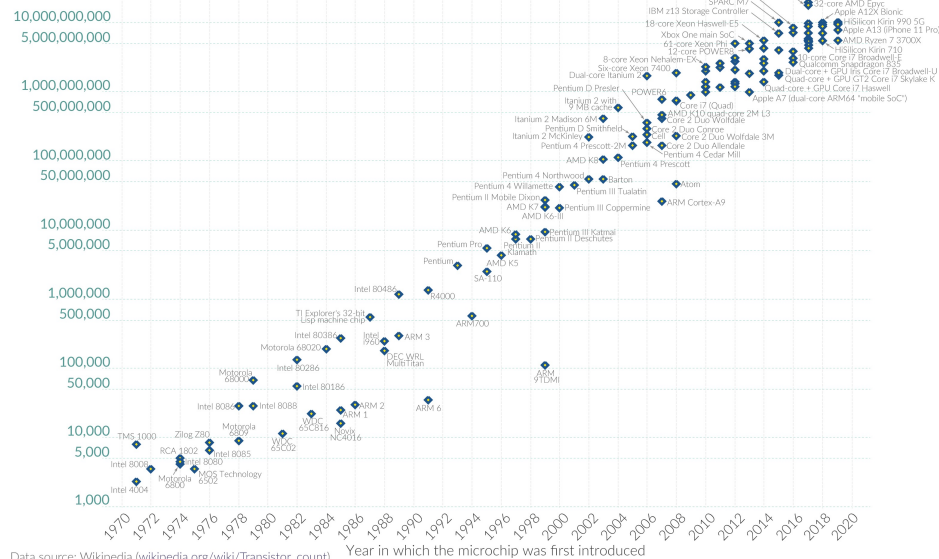


A - 18

Moore's Law: The number of transistors on microchips doubles every two years
 Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers.



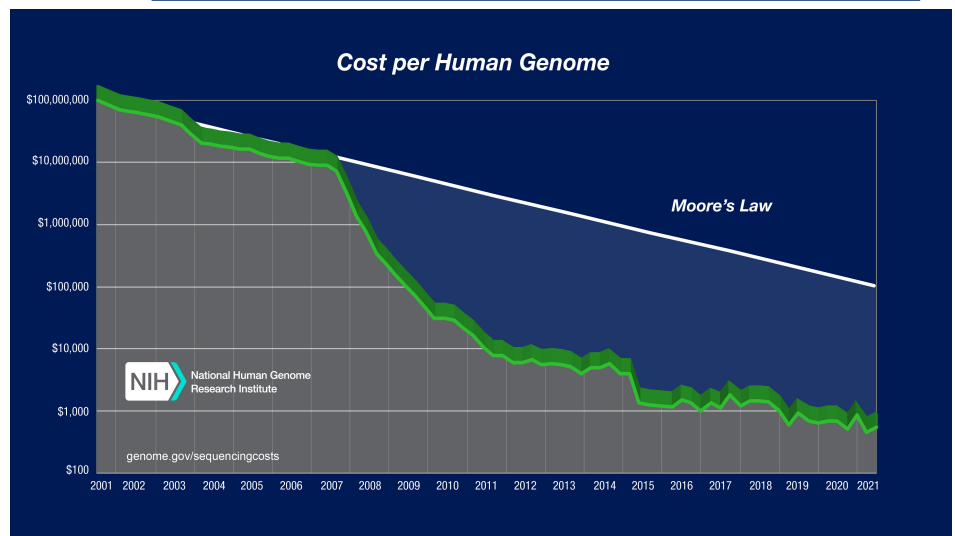
Transistor count
 50,000,000,000



Data source: Wikipedia (wikipedia.org/wiki/Transistor_count)
 OurWorldinData.org - Research and data to make progress against the world's largest problems. Licensed under CC-BY by the authors Hannah Ritchie and Max Roser.



Cost of Sequencing a Human Genome



A - 20



Recurring Themes

- Bioinformatic tools are often hypotheses-generating
- Determining statistical significance of results generated by bioinformatic tools is useful
- Properties of data guide choice of algorithm
- Some problems are solved exactly or optimally. Other problems are addressed using a *heuristic* approach.
- Many computational approaches are improved by incorporating additional biological insights into their underlying method or model.
- Recent advances have allowed scientists to gather large amounts of, often heterogeneous, data. One of the roles of bioinformatic tools is efficient analysis of large data sets with the aim of extracting new biological insights.