

CMSC423: Bioinformatic Algorithms, Databases and Tools

Dr. Todd Treangen
Lecture #2
9/4/2012

Molecular Biology primer

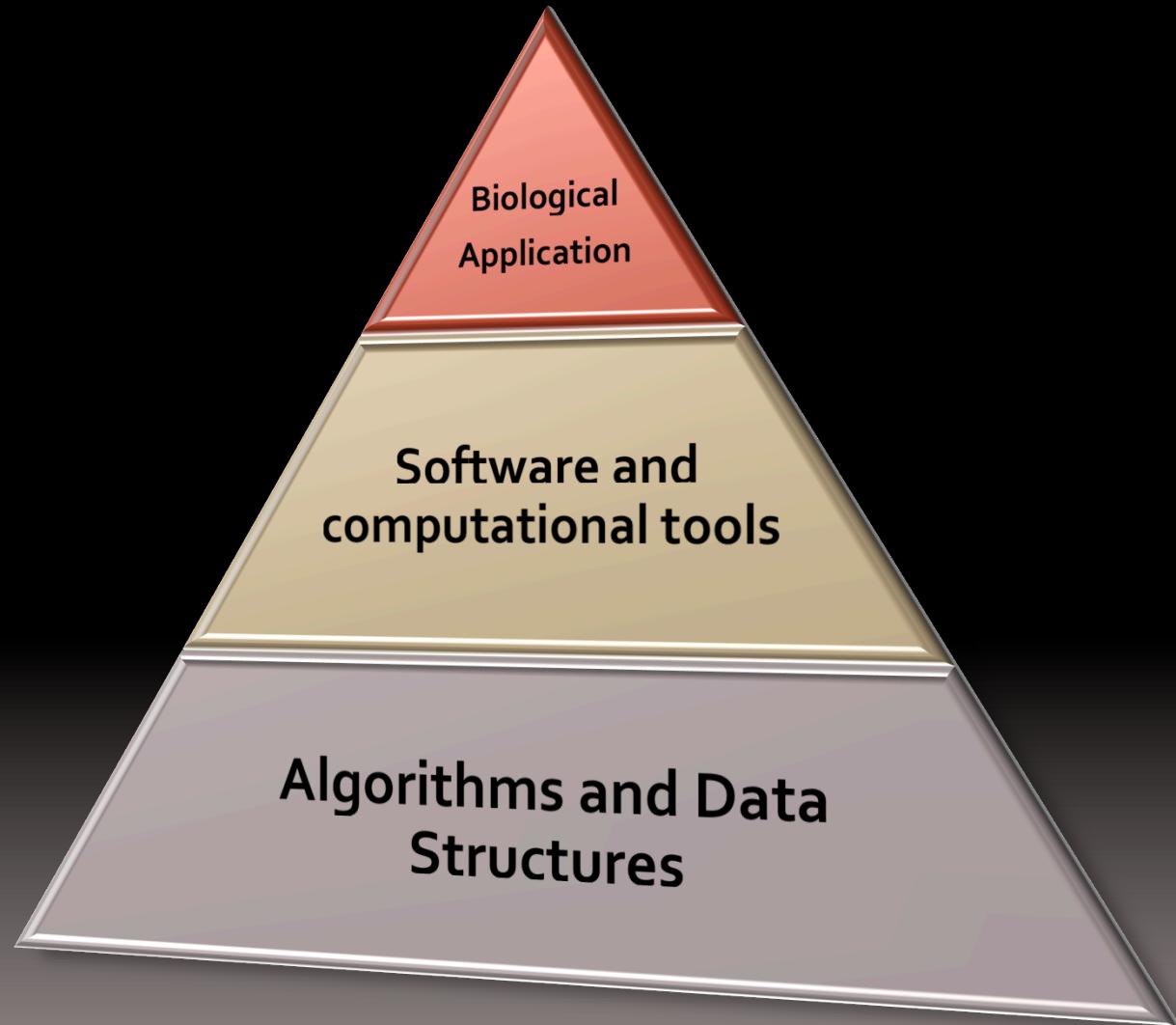
Admin...

- Review website
 - Recommend books
 - Supplemental reading material
- Have you tried your glue accounts?
 - Part of HW1 will require you to use glue system
 - Let's try it out
- Issues/concerns/questions about class and policies?
 - Attendance policy important

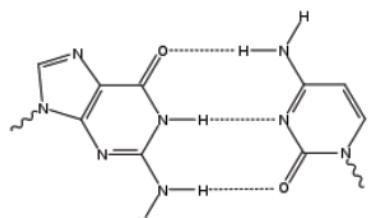
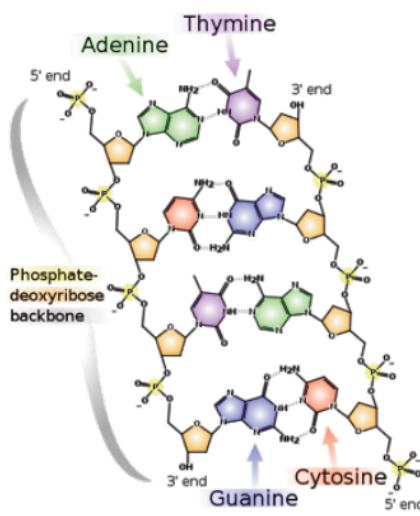
Admin continued...

- TA
 - Milad Gholami
 - Mondays & Wednesdays 9-11am
 - Everyone know where the TA office is located?
- Grades will be available through the grade server
grades.cs.umd.edu
- **Frontiers in Genomics symposium (free!)**
 - **Excellent speakers**
 - **October 18th, from 8:45am to 3pm**
 - **RSVP to igs-event@som.umd.edu**

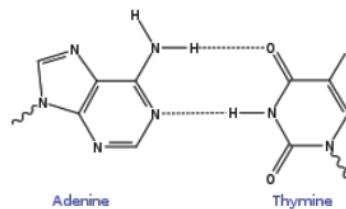
Bioinformatics



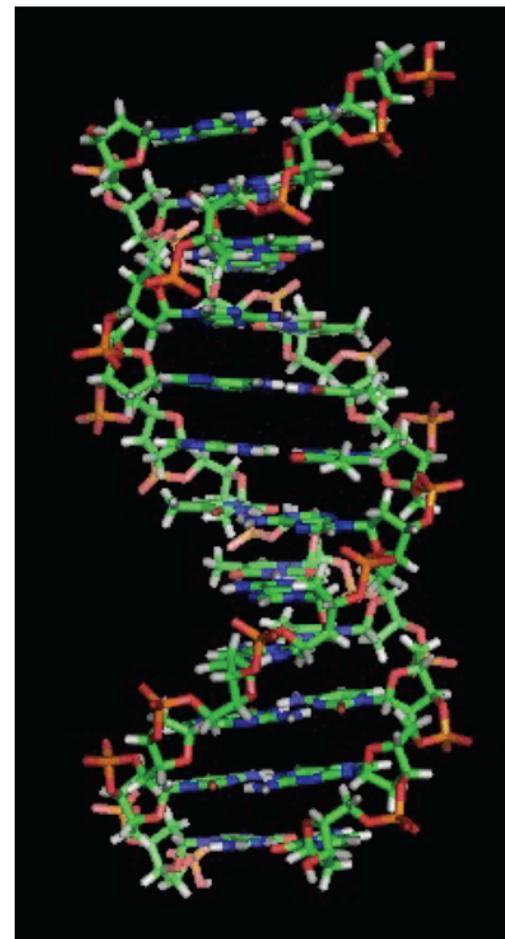
DNA



G C



A T



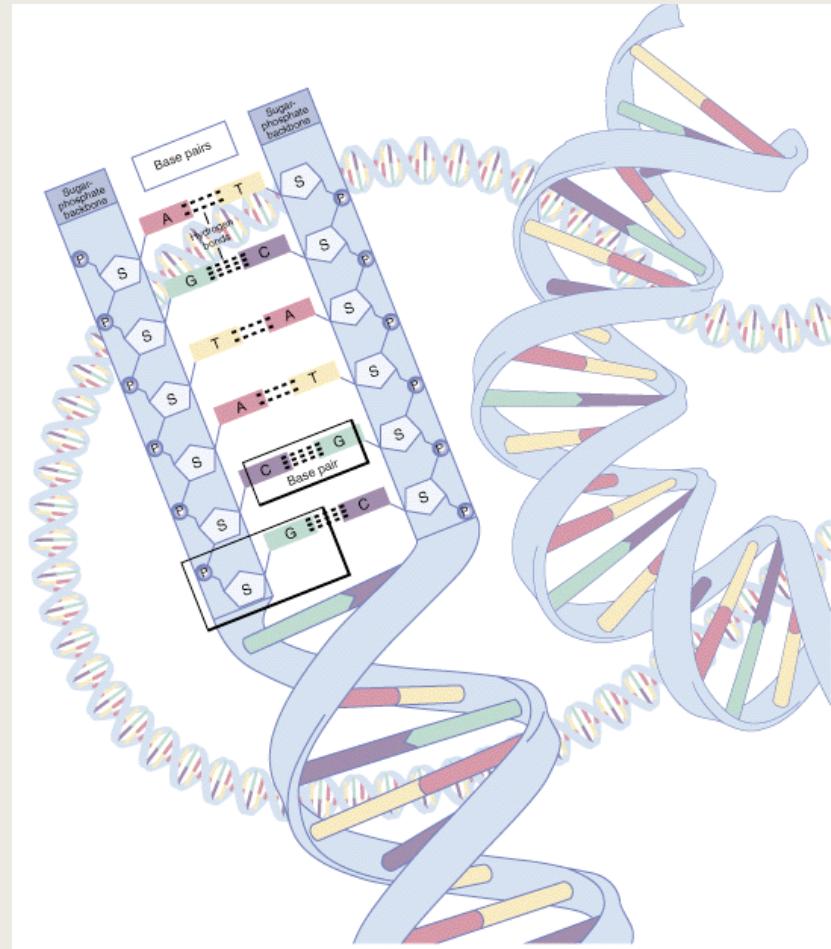
Documents of evolutionary history

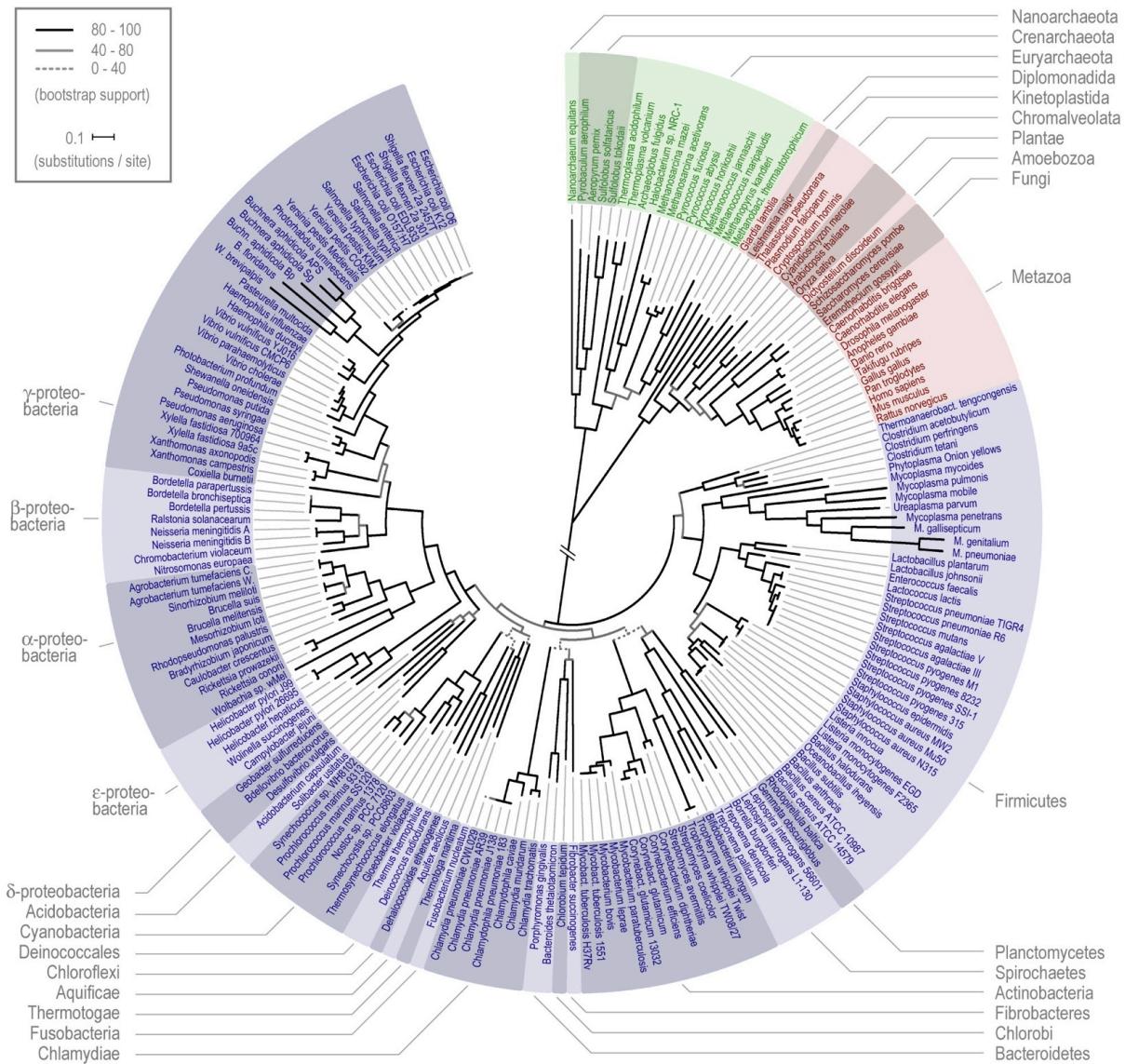
The four nucleotides in DNA contain the bases: guanine (**G**), cytosine (**C**), adenine (**A**), and uracil (**U**), read as thymine (**T**) in DNA.

The **genome** of an organism contains its hereditary information and is encoded in its DNA .

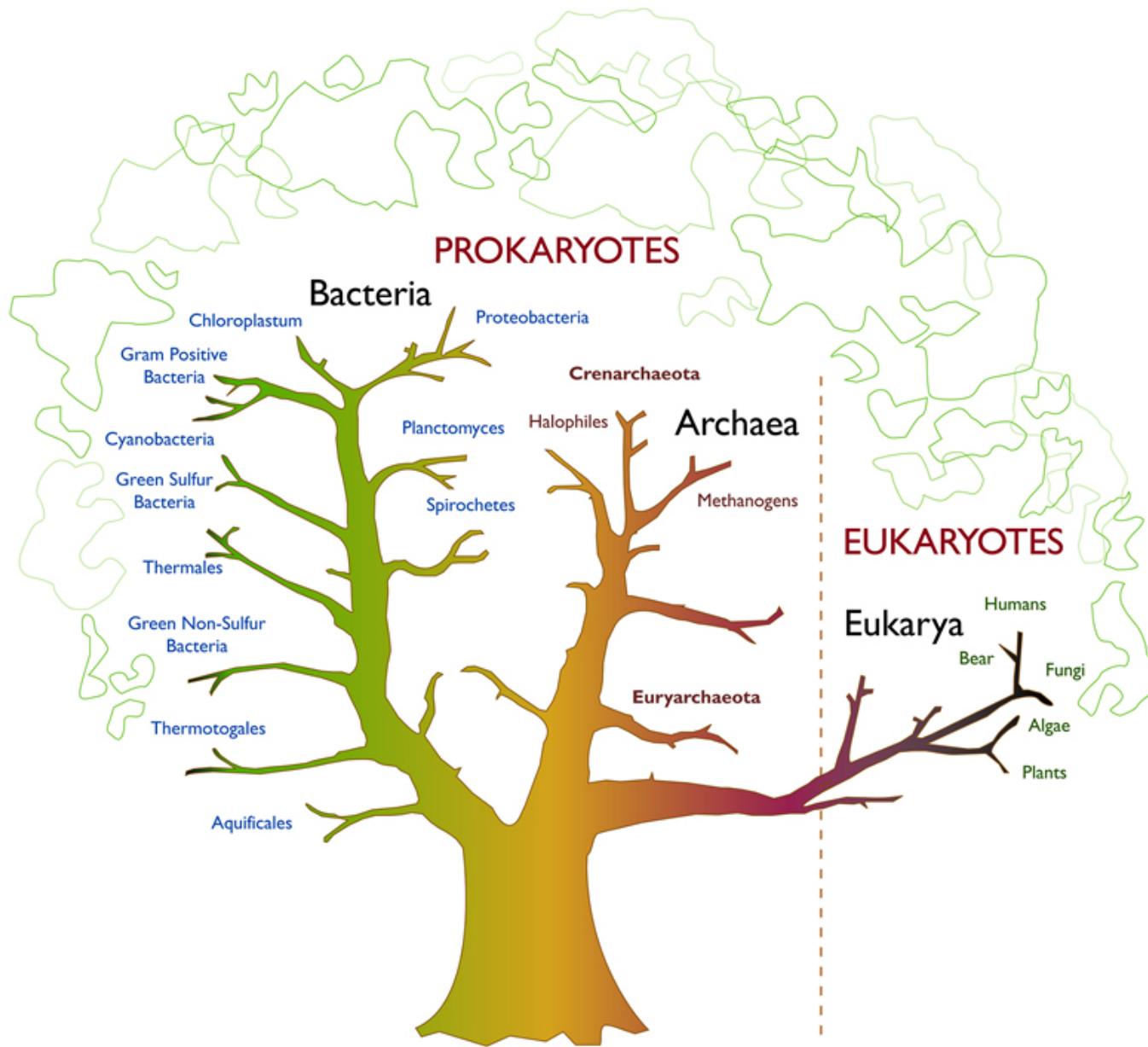
In 1965, Zuckerkandl and Pauling described an organisms DNA as "Documents of Evolutionary History"

See reading material !





Tree of Life



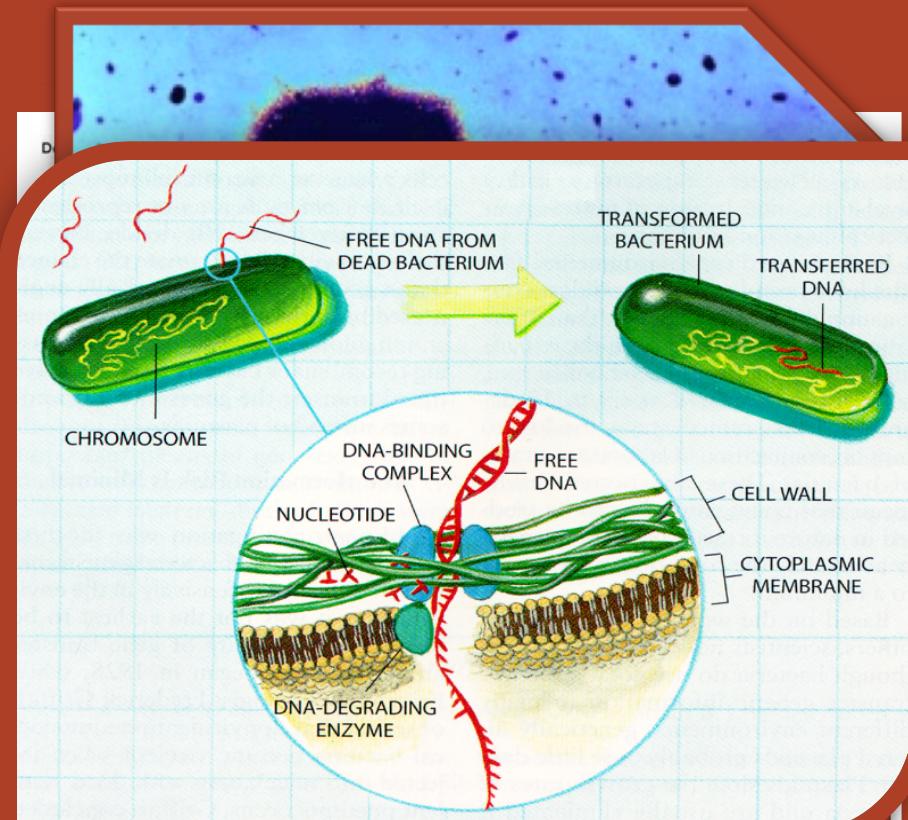
However, there are some exceptions



Genetic Exchange in bacteria

There are three main ways which a bacteria can gain new genetic material (DNA) from donor bacterium:

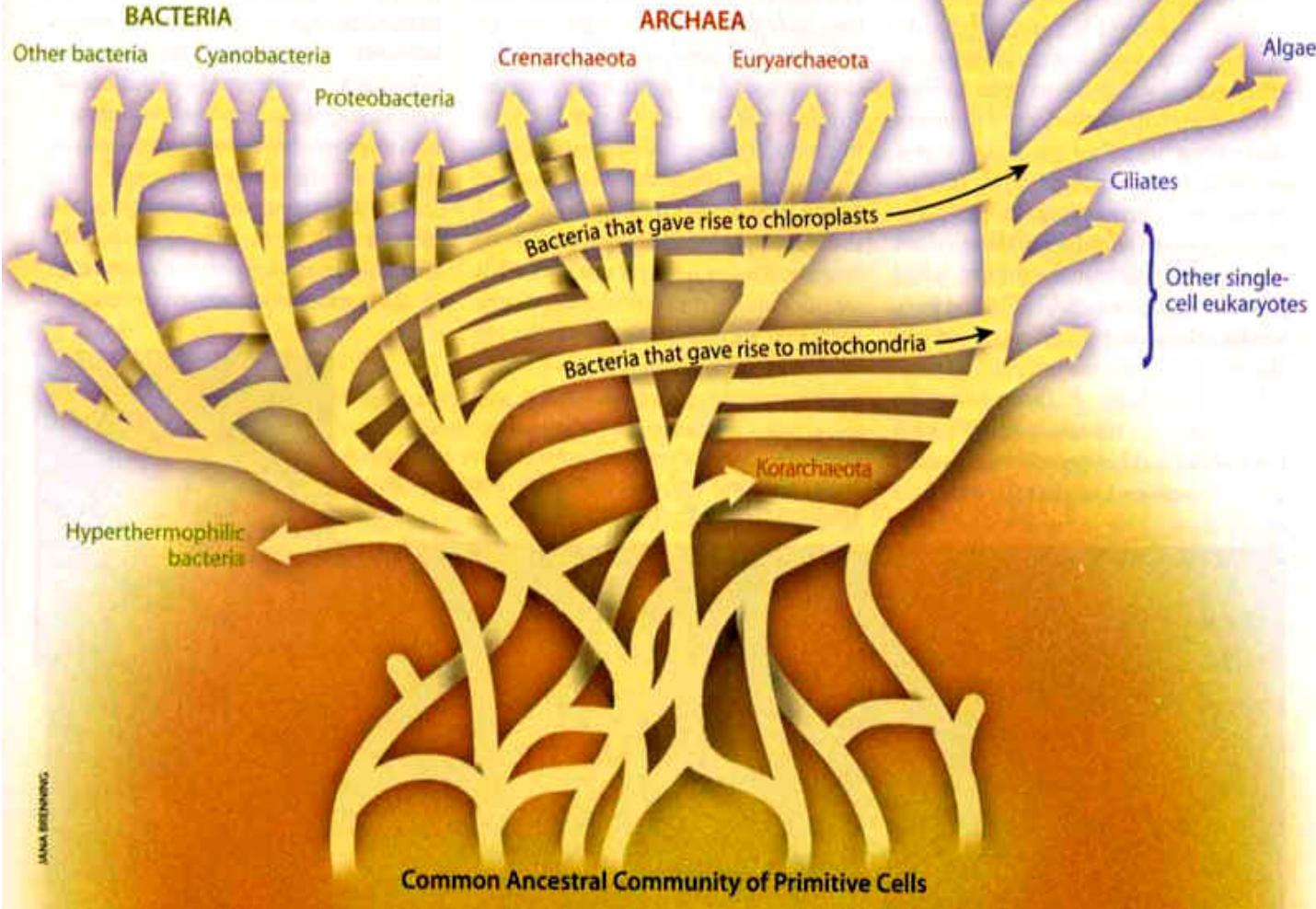
- 1. Conjugation: Transfer from one bacteria to another via a pilus. (Lederberg and Tatum, 1946)
- 2. Transduction: Transfer of DNA is mediated by a bacteriophage . (Zinder and Linderberg, 1952)
- 3. Transformation: DNA is taken up from the environment
(Griffith 1928)



Miller RV (1996) Sci. Am. 278, 66-71

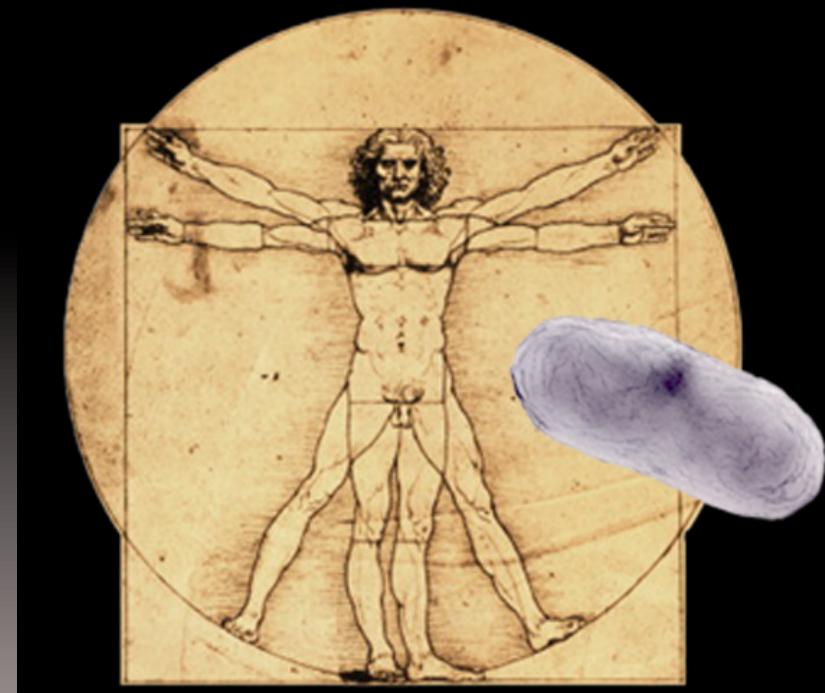
Benjamin-Cummings

REVISED "TREE" OF LIFE retains a treelike structure at the top of the eukaryotic domain and acknowledges that eukaryotes obtained mitochondria and chloroplasts from bacteria. But it also includes an extensive network of untreelike links between branches. Those links have been inserted somewhat randomly to symbolize the rampant lateral gene transfer of single or multiple genes that has always occurred between unicellular organisms. This "tree" also lacks a single cell at the root; the three major domains of life probably arose from a population of primitive cells that differed in their genes.

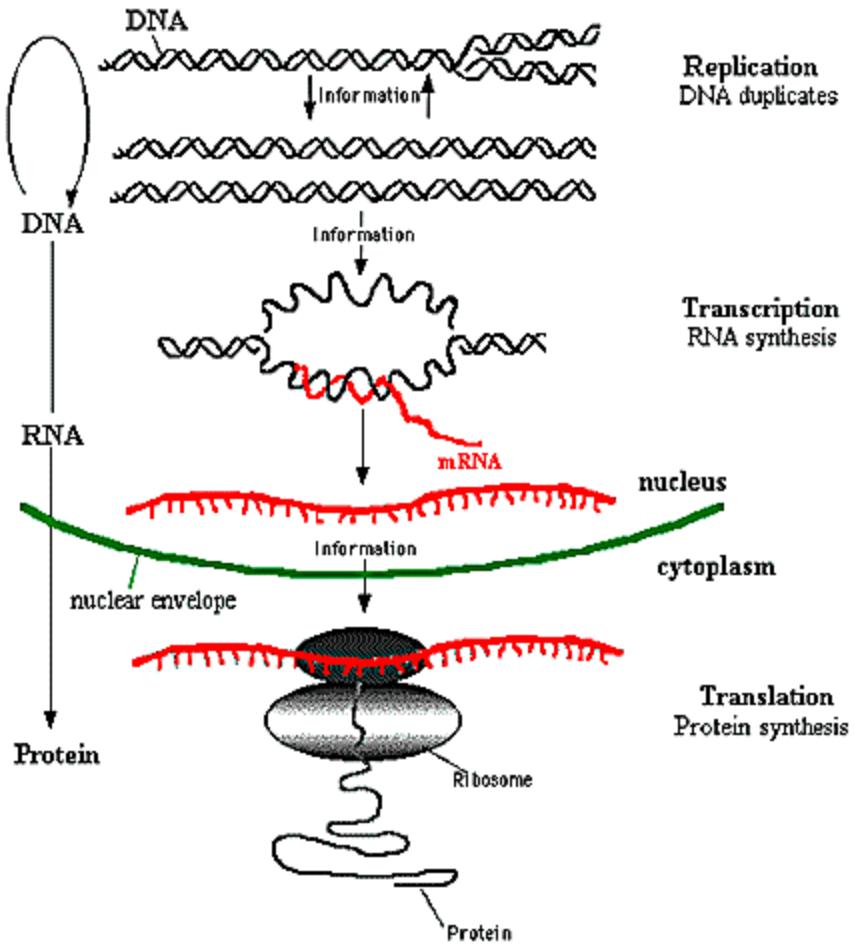


Extended view of ourselves as a lifeform

- We are composite of species: a ‘supra-organism’
- Our microbial census exceeds the total number of our own human cells by ~10 fold
- Our largest collection of microbes resides in the intestine (~10-100 trillion organisms)
- The aggregate genomes of these gut species (microbiome) may contain >100 fold more genes than our ‘own’ genome
- The microbiome is an integral part of our genetic landscape (‘human metagenome’) and of our genetic evolution

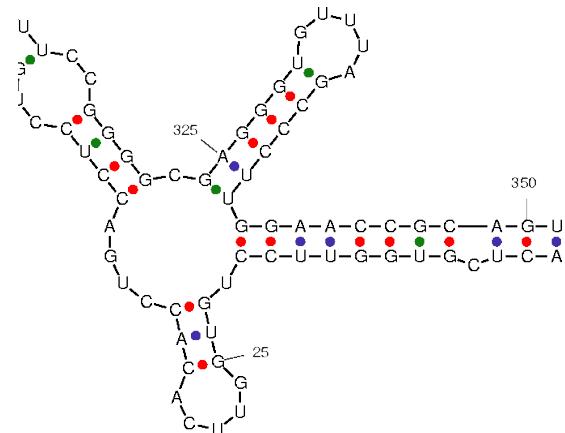


Central dogma

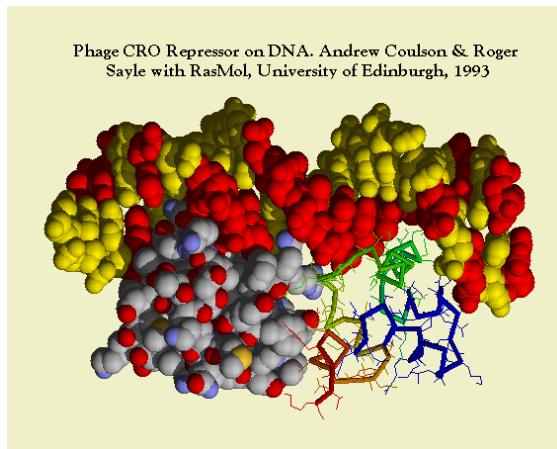


The Central Dogma of Molecular Biology

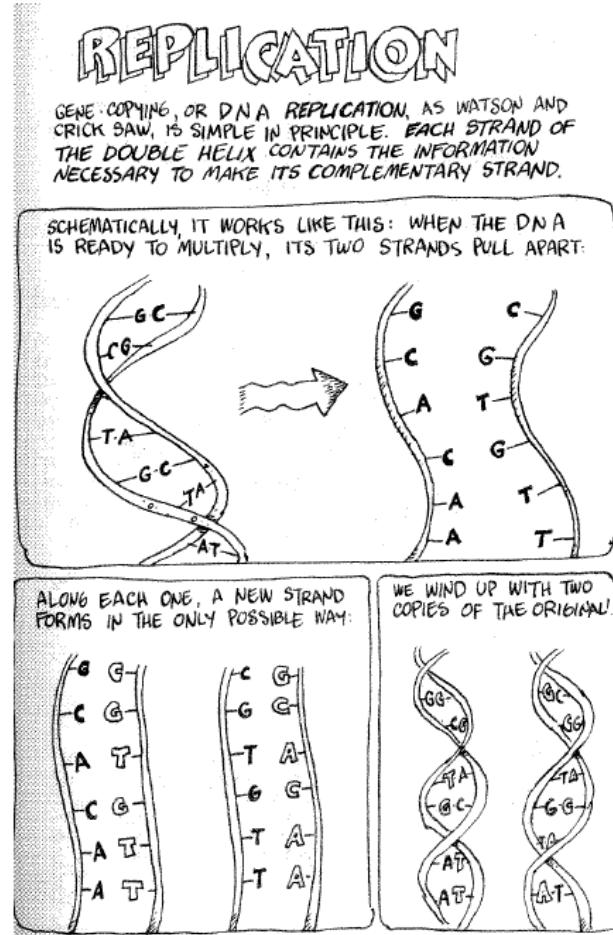
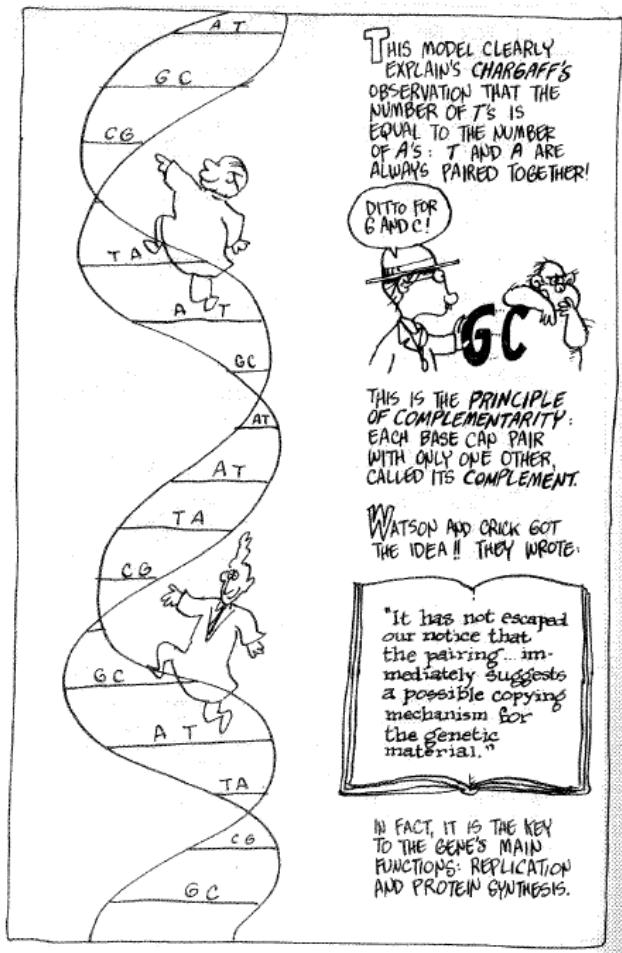
AGGTACGCGTACCTGACAGG



Phage CRO Repressor on DNA. Andrew Coulson & Roger Sayle with RasMol, University of Edinburgh, 1993



DNA Replication



The Cartoon Guide to Genetics
Larry Gonick & Mark Wheelis, 1983

Genes, transcription, translation

- DNA – RNA - Thymine replaced by Uracil (T-U)
- The transcribed segments are called genes

ACCGUACC**AUG**UUAA . . . AUAGGC**UGA**GCA

- AUG – start codon (also amino-acid Methionine)
- UAA, UAG, UGA – stop codons
- Genes are read in sets of 3 nucleotides during translation – $4^3 = 64$ possible combinations
- Each combination codes for one of 20 amino-acids – the building blocks for proteins

Amino-acid translation table

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

DNA in the computer

- FASTA/multi-FASTA file format

```
>gi|110227054|gb|AE004091.2| Pseudomonas aeruginosa PA01, complete genome  
TTTAAAGAGACCGGCGATTCTAGTGAAATCGAACGGCAGGTCAATTCCAACCAGCGATGACGTAATAG  
ATAGATAACAAGGAAGTCATTCTAAAGGATAGAAACGGTTAATGCTCTGGGACGGCGCTTTCT  
GTGCATAACTCGATGAAGCCCAGCAATTGCGTCTCCGGCAGGCAAAAGGTTGTCGAGAACCGGTGT  
CGAGGCTGTTCCCTGAGCGAAGCCTGGGATGAACGAGATGGTTATCCACAGCGGTTTTCCACA  
CGGCTGTGCGAGGGATGTACCCCTCAAAGCAAGGTTATCCACAAAGTCCAGGACGACCGTCCGTCG
```

- Parsers easy to write, also available in a variety of software libraries

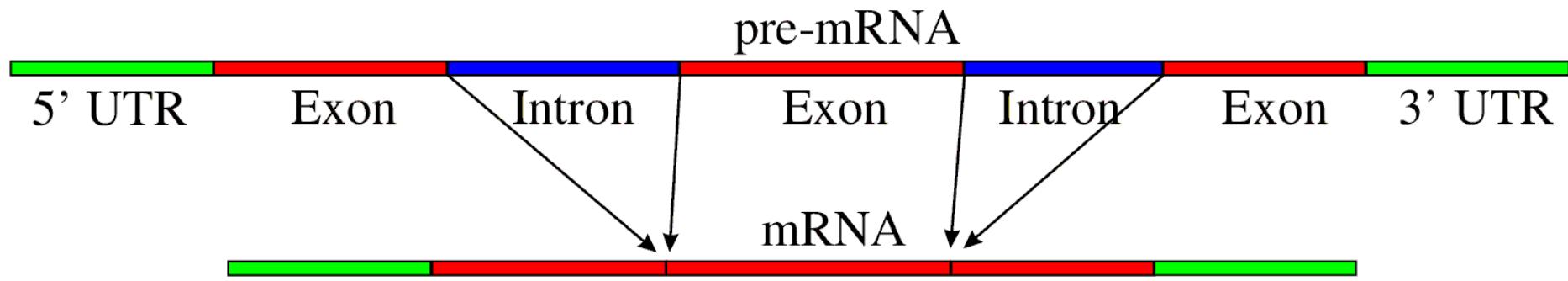
Genes/proteins in the computer

- >gi|15596155|ref|NP_249649.1| basic amino acid,
▪ MKVMKWSAIALAVSAGSTQFAVADAFVSDQAEAKGFIEDSSL DLLRNYYFNRDGKSGSGDRVDWTQGFL
▪ TTYESGFTQGTVGFGVDAFGYLGLKLDGTSDKTGTGNLPVMNDGKPRDDYSRAGGAVKVRISKMLKWGE
▪ MQPTAPVFAAGGSRLFPQTATGFQLQSSEFEGLDLEAGHFTEGKEPTTVKSRGELYATYAGETAKSADFI
▪ GGRYAITDNL SASLYGAELEDIYRQYYLNSNYTIPLASDQSLGFDFNIYRTNDEGKAKAGDISNTTWSLA
▪ AAYTLDAHTFTLAYQKVHGDQPFDYIGFGRNGSGAGGDSIFLANSVQYSDFNGPGEKSWQARYDLNLASY
▪ GVPGLTFMVR YINGKDIDGTKMSDNNVGYKNYGYGEDGKHETNLEAKYVVQSGPAKDL SFRIRQAWHRA
▪ NADQGEGDQNEFRLIVDYPLSIL
- Same FASTA/multi-FASTA but with bigger alphabet

Genes/proteins in the computer

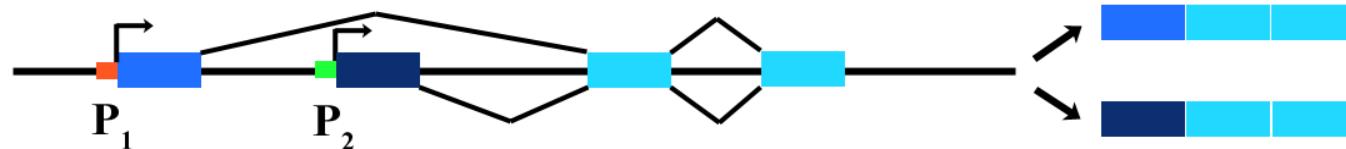
- gene complement(1043983..1045314)
▪ /gene="oprD"
▪ /locus_tag="PA0958"
- CDS complement(1043983..1045314)
▪ /gene="oprD"
▪ /locus_tag="PA0958"
▪ /note="Product name confidence: class 1 (Function experimentally demonstrated in *P. aeruginosa*)"
▪ /codon_start=1
▪ /transl_table=11
▪ /product="Basic amino acid, basic peptide and imipenem outer membrane porin OprD precursor"
▪ /protein_id="AAG04347.1"
▪ /db_xref="GI:9946864"
- GenBank file format

Translation – complications

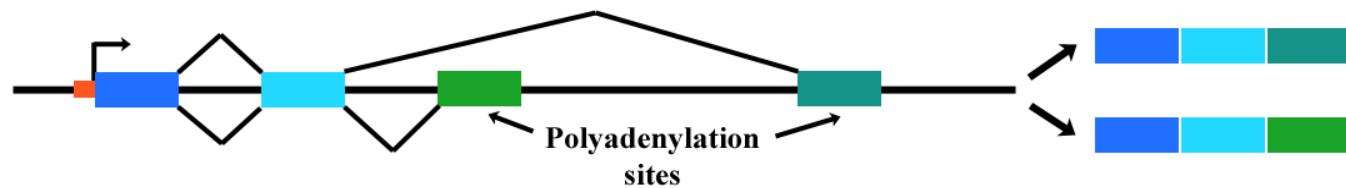


Alternative splicing examples

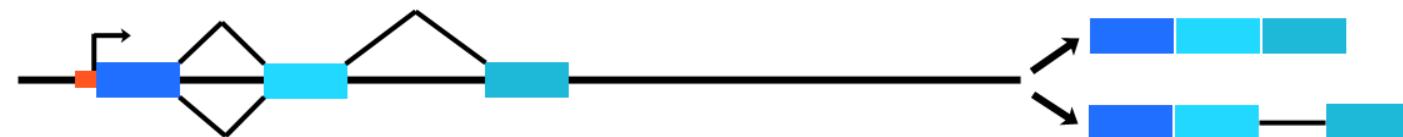
(a) Alternative selection of promoters (e.g., *myosin* primary transcript)



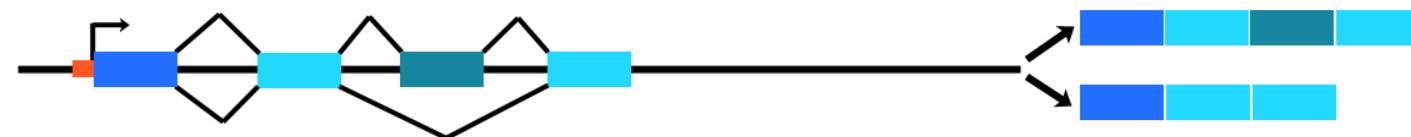
(b) Alternative selection of cleavage/polyadenylation sites (e.g., *tropomyosin* transcript)



(c) Intron retaining mode (e.g., *transposase* primary transcript)



(d) Exon cassette mode (e.g., *troponin* primary transcript)



RECAP

- DNA is a string formed with letters A, C, T, G (called nucleotides or bases)
- DNA is double-stranded – allows replication: transfer of genetic “code” from parents to offspring
- DNA is naturally oriented from 5' to 3' and the two strands are anti-parallel
- If you know the sequence of one strand, you can obtain the sequence of the other by reverse-complementation

5' AGACCTAGTGCACGGCTACTACC 3'

5' CCATCATCGGCACGTGATCCAGA 3' Reverse

5' GGTAGTAGCCGTGCACTAGGTCT 3' Complement

RECAP

- Central Dogma of molecular biology:
 - DNA – RNA (transcription)
 - RNA – Protein (translation)
- The transcribed segments of DNA are called “genes”
- Translation occurs in sets of 3 nucleotides – codons
- Each codon encodes one of 20 amino-acids and 3 stop-codons
- In eukaryotes the genes may be split into multiple exons, separated by introns: DNA segments that will not get translated
- The protein is translated from an RNA representing the concatenation of the exons of the gene₂₃

The “new” biology

- DNA is not the only heritable information
 - Epigenetic information: RNA molecules, DNA methylation patterns (affects coiling on DNA on histones)
- Complex regulation patterns
 - Genes turn on other genes
 - Genes inhibit other genes
 - RNA interference – small RNA molecules can destroy specific transcripts (down-regulate production)

Playing with DNA

- Biologists can:
- Cut the DNA – restriction enzymes (often palindromes)
(Nobel prize – Arber, Nathans, Smith)

5'GAATTC
3'CTTAAG

5'---G AATTC---3'
3'---CTTAA G---5'

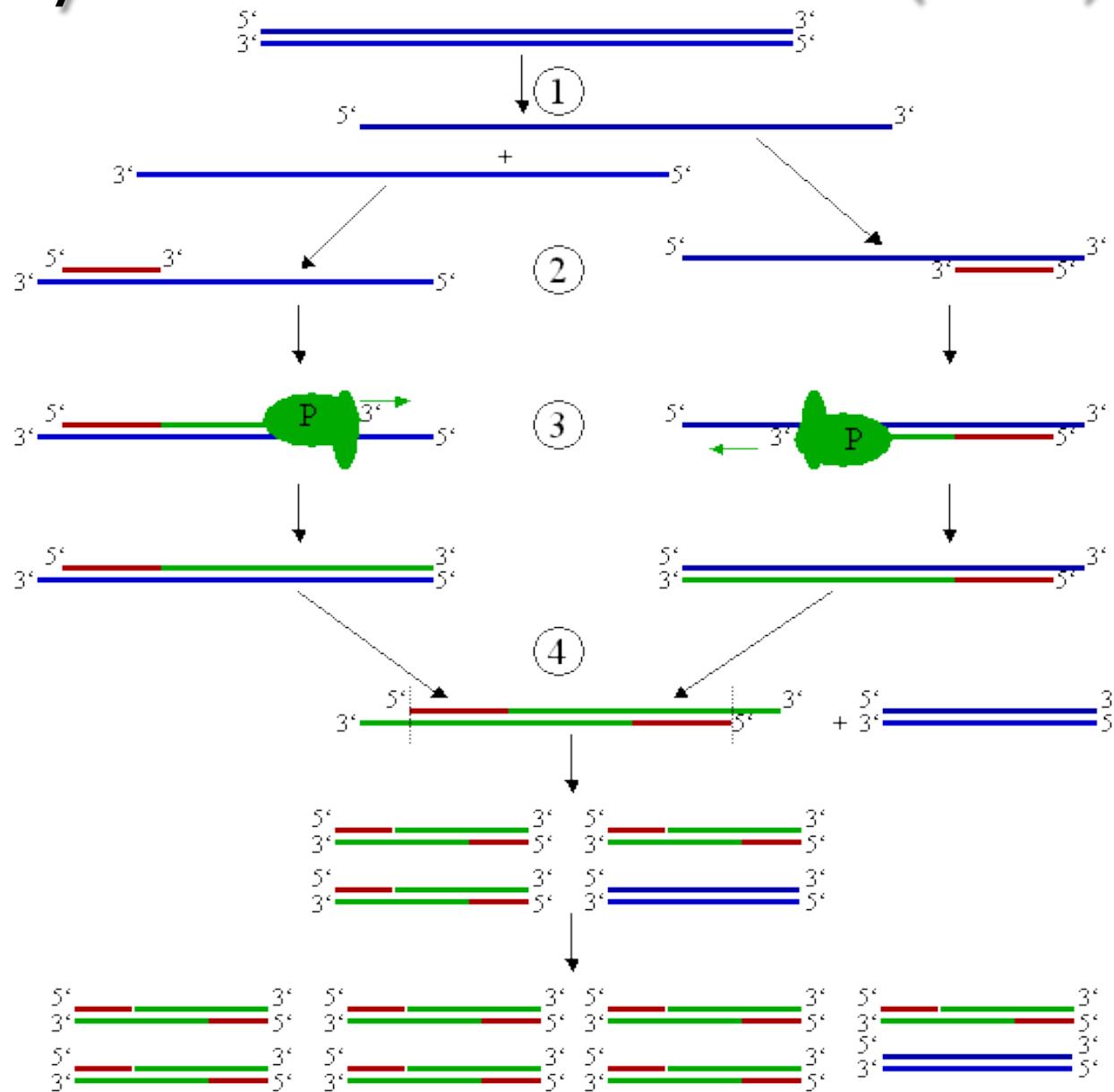
- Attach “things” to DNA (either single or double-strand)

TAGGCACGTTGCAACTACGGC

TGCAACGT

- “Amplify” DNA – Polymerase Chain Reaction
(Nobel prize – Mullis)

Polymerase chain reaction (PCR)



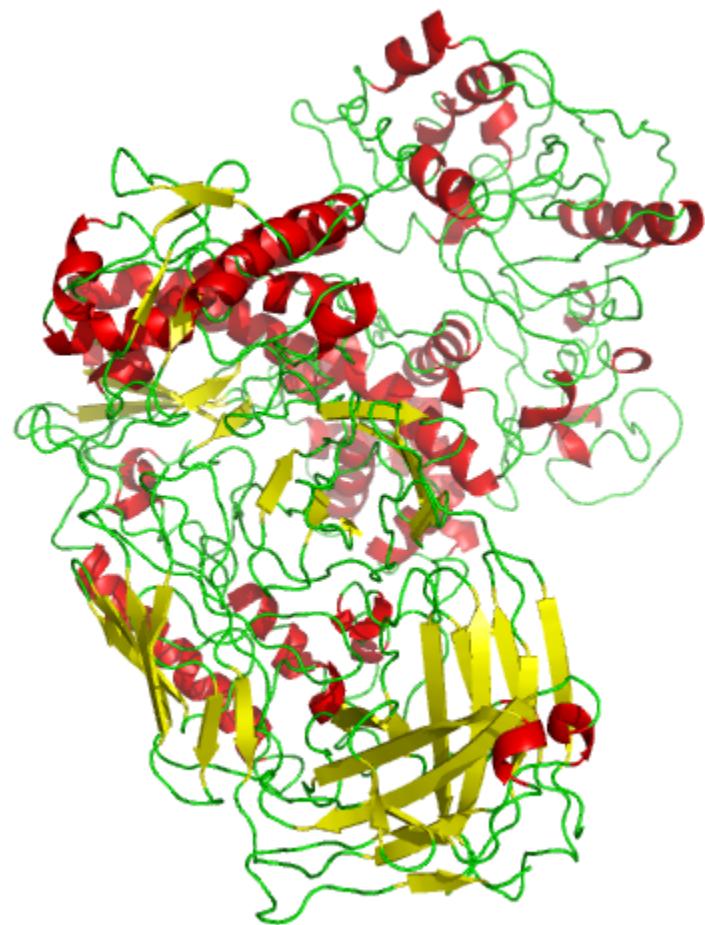
1. Denature

2. Anneal (attach primer)

3. Extend

4. Repeat

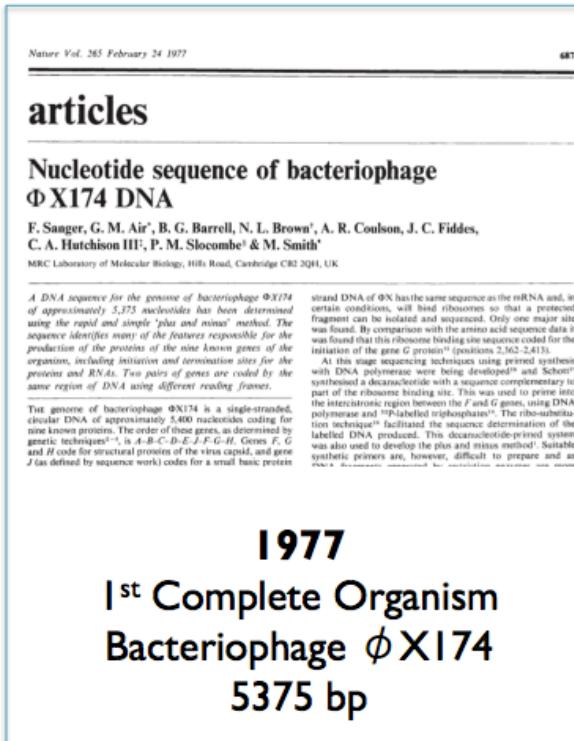
Taq polimerase



DNA sequencing

- Most techniques “trick” the polymerase into revealing the sequence
- The traditional method – Sanger sequencing – based on “terminator” bases – prevent the polymerase from extending the DNA
- Sanger sequencing is essentially PCR + terminator bases
- Other methods “spy” on the polymerase as it incorporates nucleotides

Milestones in Molecular Biology



Radioactive Chain Termination
5000bp / week / person

<http://en.wikipedia.org/wiki/File:Sequencing.jpg>
<http://www.answers.com/topic/automated-sequencer>

Nucleotide sequence of bacteriophage ϕ X174 DNA
Sanger, F. et al. (1977) *Nature*. 265: 687 - 695

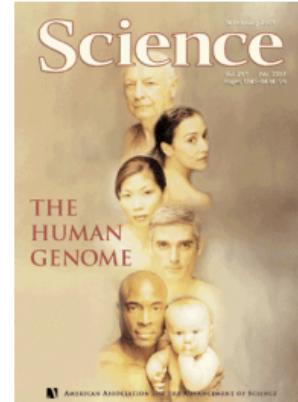
Milestones in Molecular Biology



1995
Fleischmann et al.
1st Free Living Organism
TIGR Assembler. 1.8Mbp



2000
Myers et al.
1st Large WGS Assembly.
Celera Assembler. 116 Mbp



2001
Venter et al. / IHGSC
Human Genome
Celera Assembler. 2.9 Gbp

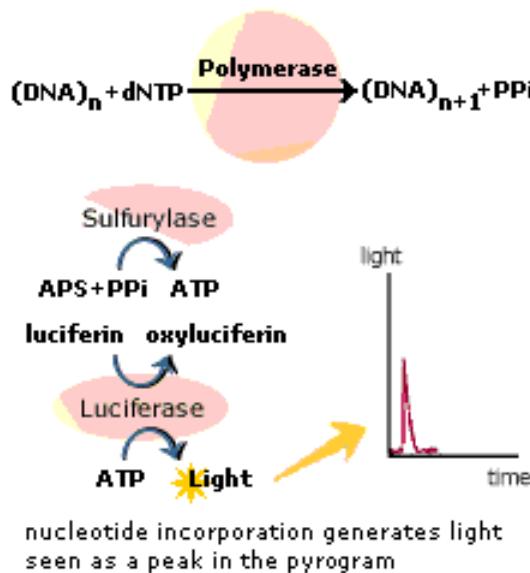
ABI 3700: 500 bp reads x 768 samples / day = 384,000 bp / day.

"The machine was so revolutionary that it could decode in a single day the same amount of genetic material that most DNA labs could produce in a year." J. Craig Venter

The future of sequencing

- Single molecule sequencing - current technology requires many copies of DNA being sequenced - requires DNA amplification
- Massively-parallel sequencing - 100k sequencing reactions occurring at the same time

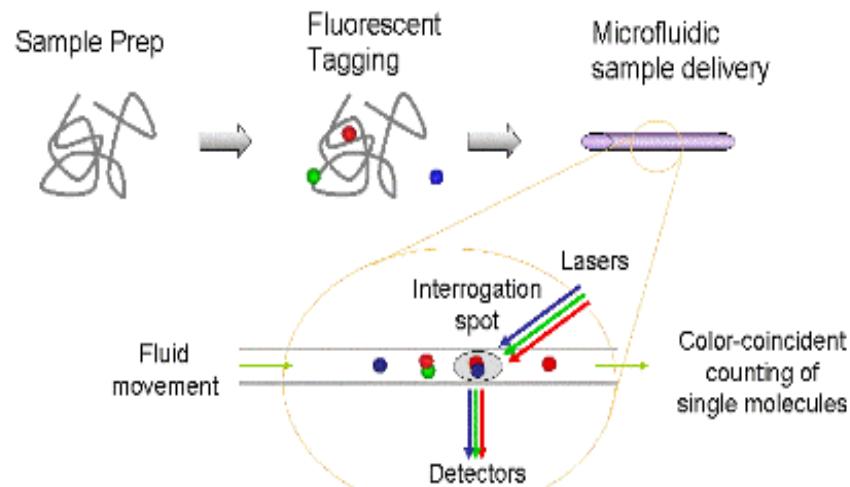
Sequencing by synthesis



TCTAAT^AG^A

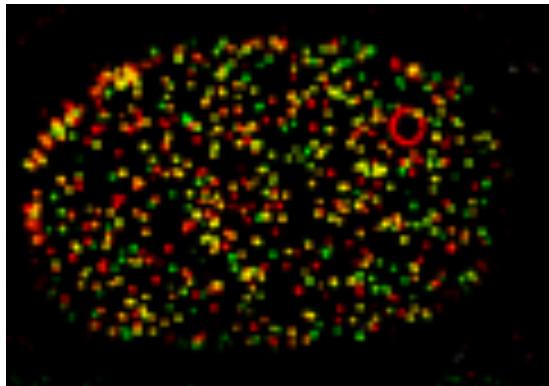
AGATTATCTAACAGCTACCCTTCCATCA

Micro-fluidics



The future of sequencing

Massively parallel sequencing



<http://arep.med.harvard.edu/>

- each spot is a molecule or amplified from one molecule
- image processing used to track molecules during sequencing by synthesis
- often micro-fluidics/lab-on-a-chip used

- More on this in two weeks!

The evolution of DNA sequencing

Since	Technology	Read length	Throughput/run	Throughput/hour	cost/run
1977-	Sanger sequencing	> 1000bp	4hr 400-500 kbp	100 kbp	\$200
2005-	454 pyrosequencing	250-400bp	4hr 100-500 Mbp	25-100 Mbp	\$13,000
2006-	Illumina/Solexa	50-100bp	3 days 2-3 Gbp	25-40 Mbp	\$3,000
2007-	ABI SOLiD	35-50bp	3 days 6-20 Gbp	75-250 Mbp	est. \$3-5,000
2010-	Illumina HiSeq 2000	2X100 bp	11 days 600 Gbp	2200 Mbp	~\$30,000
2010-	Pacific Biosciences single molecule	1-10 kbp	1 day 2.4 Gbp	100 Mbp	\$1,000
2012-	Oxford Nanopore	80-100 kbp	?	?	?

In the news (more info on website)

Published Online August 30 2012

Science DOI: 10.1126/science.1224344

< Science Express Index

 Read Full Text to Comment (0)

RESEARCH ARTICLE

A High-Coverage Genome Sequence from an Archaic Denisovan Individual

Matthias Meyer^{1,*†}, Martin Kircher^{1,*†}, Marie-Theres Gansauge¹, Heng Li², Fernando Racimo¹, Swapan Mallick^{2,3}, Joshua G. Schraiber⁴, Flora Jay⁴, Kay Prüfer¹, Cesare de Filippo¹, Peter H. Sudmant⁶, Can Alkan^{4,5}, Qiaomei Fu^{1,7}, Ron Do², Nadin Rohland^{2,3}, Arti Tandon^{2,3}, Michael Siebauer¹, Richard E. Green⁸, Katarzyna Bryc³, Adrian W. Briggs³, Udo Stenzel¹, Jesse Dabney¹, Jay Shendure⁶, Jacob Kitzman⁶, Michael F. Hammer⁹, Michael V. Shunkov¹⁰, Anatoli P. Derevianko¹⁰, Nick Patterson², Aida M. Andrés¹, Evan E. Eichler^{6,11}, Montgomery Slatkin⁴, David Reich^{2,3,*†}, Janet Kelso¹, Svante Pääbo^{1,‡}

 Author Affiliations

 Author Notes

 To whom correspondence should be addressed. E-mail: mmeyer@eva.mpg.de (M.M.); reich@genetics.med.harvard.edu (D.R.); paabo@eva.mpg.de (S.P.)

 These authors contributed equally to this work.

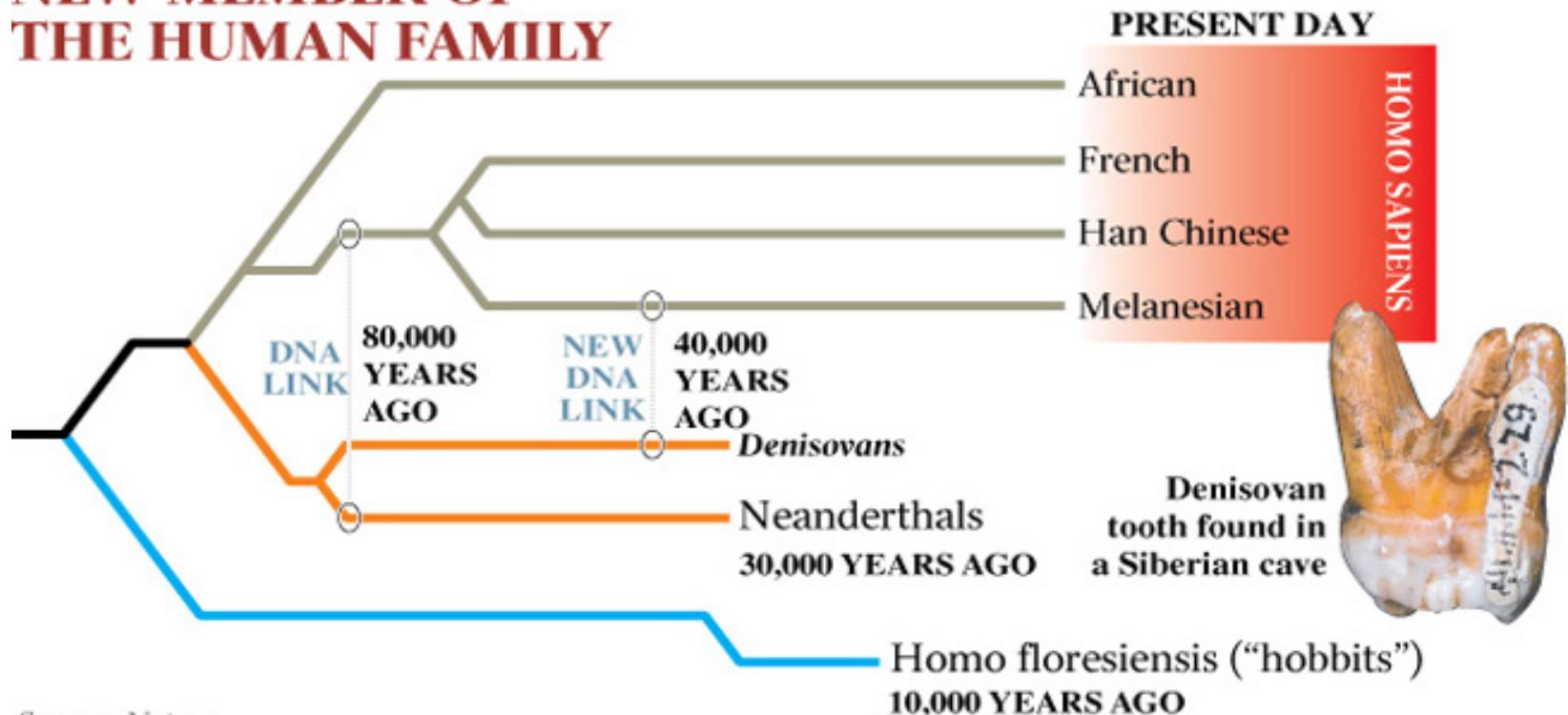
ABSTRACT

We present a DNA library preparation method that has allowed us to reconstruct a high-coverage (30X) genome sequence of a Denisovan, an extinct relative of Neandertals. The quality of this genome allows a direct estimation of Denisovan heterozygosity, indicating that genetic diversity in these archaic hominins was extremely low. It also allows tentative dating of the specimen on the basis of "missing evolution" in its genome, detailed measurements of Denisovan and Neandertal admixture into present-day human populations, and the generation of a near-complete catalog of genetic changes that swept to high frequency in modern humans since their divergence from Denisovans.

DNA sequenced from bone is 30,000 to 50,000 years old!

Documents of evolutionary history

NEW MEMBER OF THE HUMAN FAMILY



Recap

- Central dogma of biology: DNA -> RNA -> Proteins
 - DNA encodes genes, most of which encode for proteins (via the genetic code)
 - Proteins perform much of the work of the cell.
 - RNA acts as an intermediate step
(it also has other functions as well)
- Huge amount of data now available, need algorithms to make sense of it.

Notes & next lecture

- Homework #1 is posted
- Slides will be posted shortly
- Reading material is posted

- Next lecture:
 - Bioinformatics programming, DBs
 - Go over project specification
 - Might start string comparison