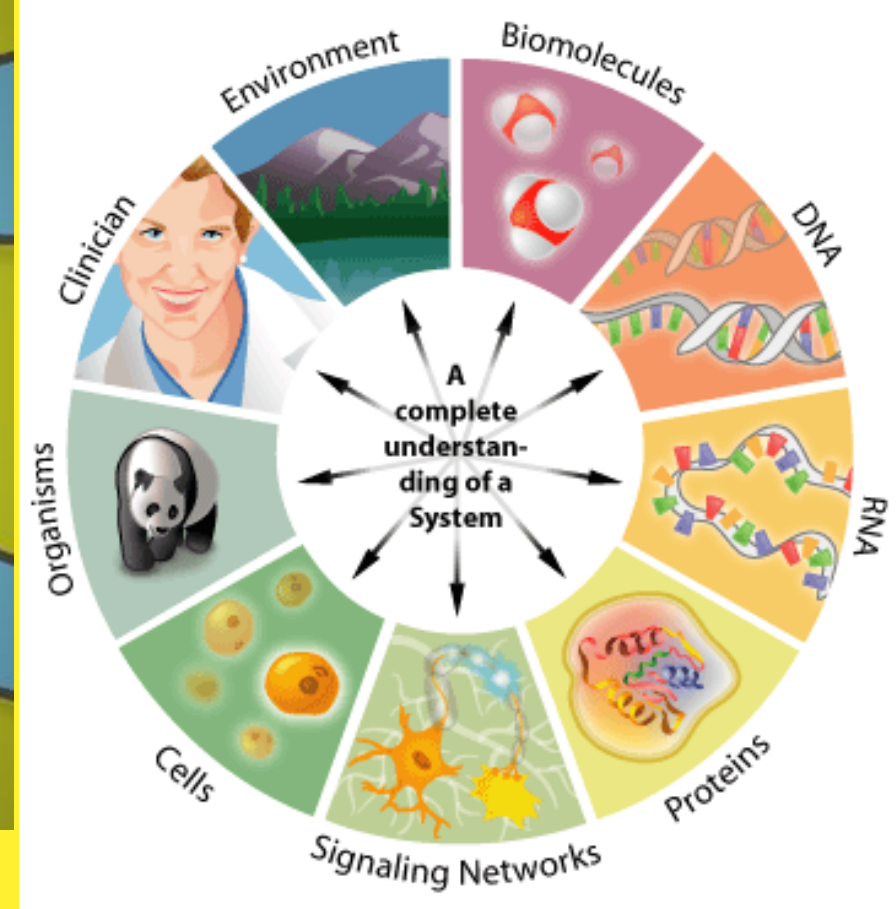
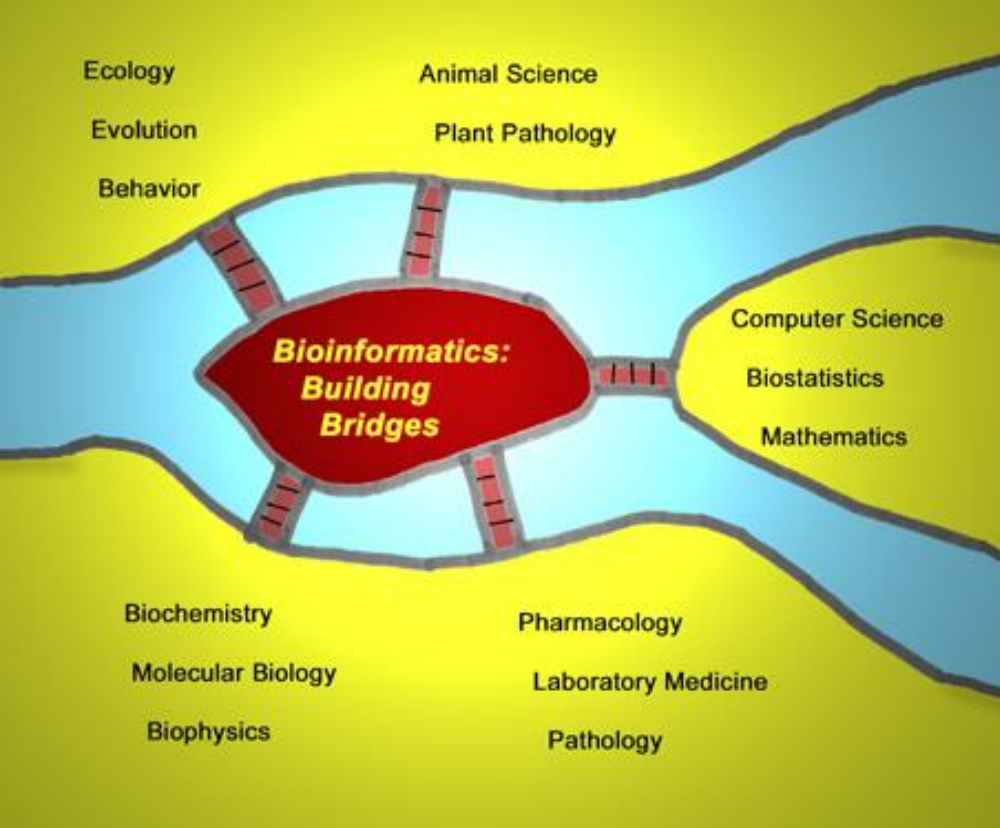


Introduction to Bioinformatics



What is Bioinformatics?

NIH – definitions

What is Bioinformatics? - ~~Research, development,~~
~~and application of computational tools and~~ on molecular
~~approaches for expanding the use of biological,~~
~~medical, behavioral, and health data, including the~~
means to acquire, ~~store, organize, archive,~~ analyze,
or visualize such data.

~~**What is Computational Biology?** The
development and application of analytical and
theoretical methods, mathematical modeling and
computational simulation techniques to the study of
biological, behavioral, and social data.~~

Molecular Bioinformatics

Molecular Bioinformatics involves the use of computational tools to discover new information in complex data sets (from the **one-dimensional** information of DNA through the **two-dimensional** information of RNA and the **three-dimensional** information of proteins, to the **four-dimensional** information of evolving living systems). **sometimes**

Bioinformatics (Oxford English Dictionary):

The branch of science concerned with information and information flow in biological systems, esp. the use of computational methods in genetics and genomics.

The field of science in which **biology**, **computer science** and **information technology** merge into a single discipline

Biologists

collect molecular data:
DNA & Protein sequences,
gene expression, etc.

Bioinformaticians

Study biological questions by
analyzing molecular data

Computer scientists

(+Mathematicians, Statisticians, etc.)
Develop tools, softwares, algorithms
to store and analyze the data.

Some biological background....

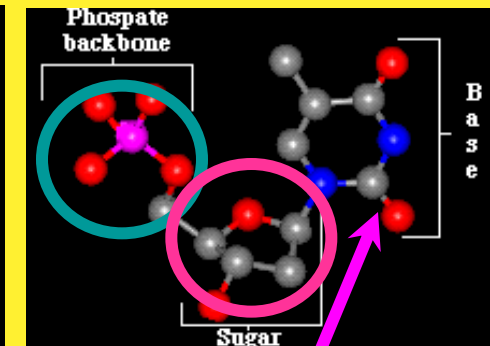
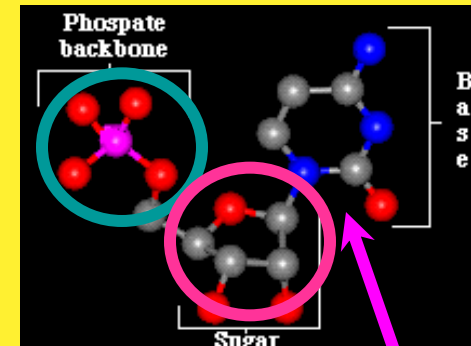
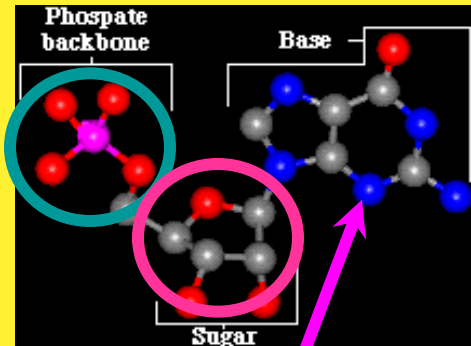
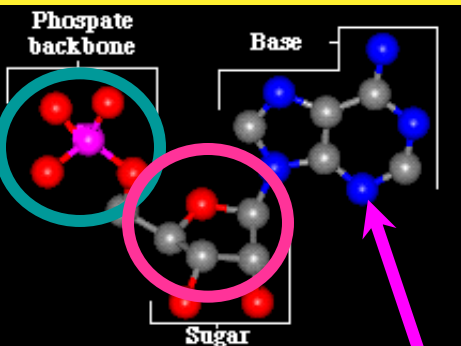
A biologist



The hereditary information of all living organisms, with the exception of some viruses, is carried by **deoxyribonucleic acid (DNA)** molecules.

2 purines:

2 pyrimidines:



adenine
(A)

guanine (G)

cytosine
(C)

thymine (T)

*two
rings*

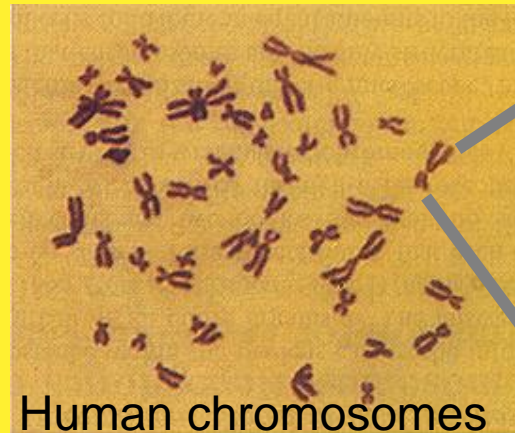
*one
ring*

The entire complement of genetic material carried by an individual is called the **genome**.

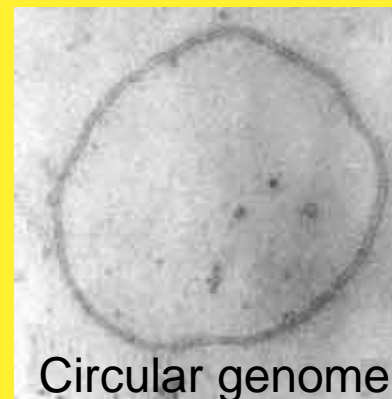
Eukaryotes may have up to 3 subcellular genomes:

1. Nuclear
2. Mitochondrial
3. Plastid

Bacteria have either circular or linear genomes and may also carry plasmids



Human chromosomes

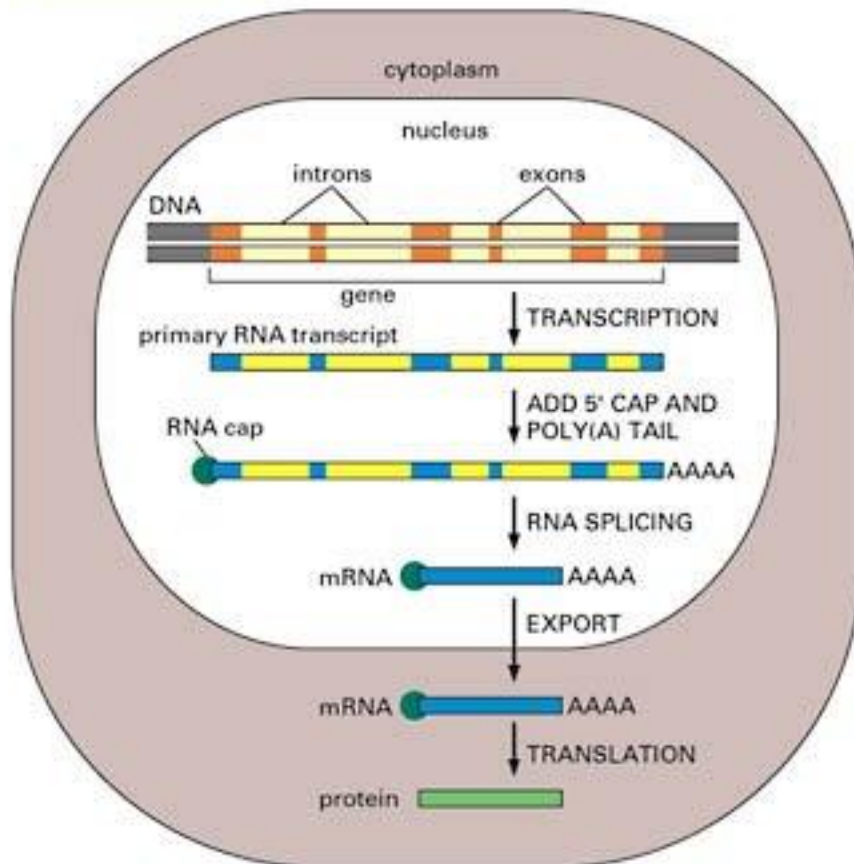


Circular genome

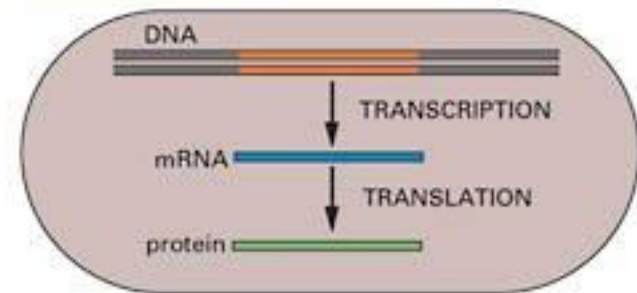
Central dogma: DNA makes RNA makes Protein

Modified dogma: DNA makes DNA and RNA, RNA makes DNA, RNA an Protein

(A) EUCARYOTES

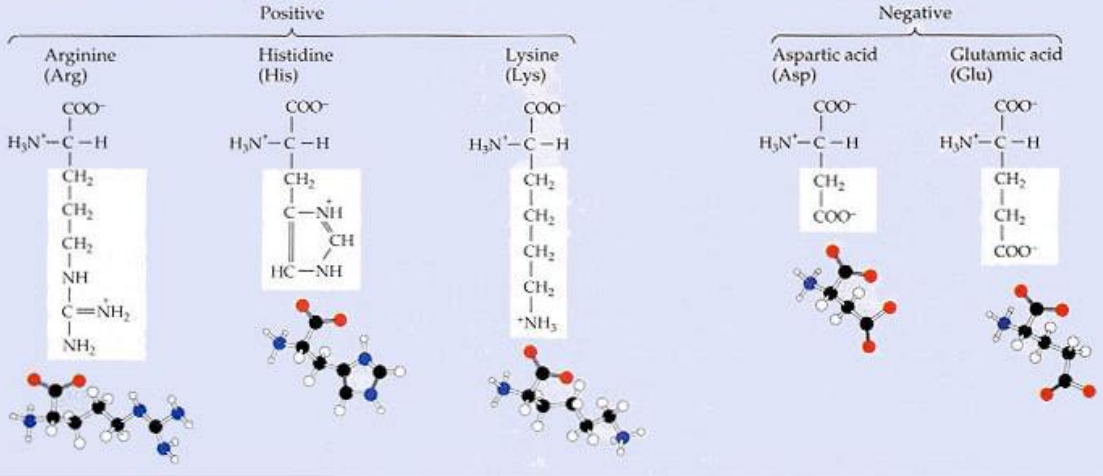


(B) PROCARYOTES

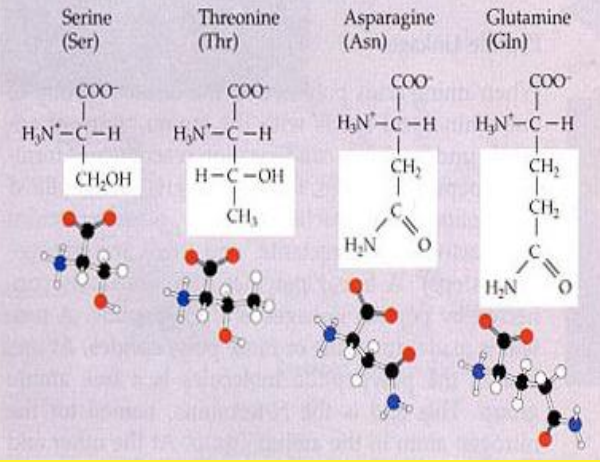


Amino acids - The protein building blocks

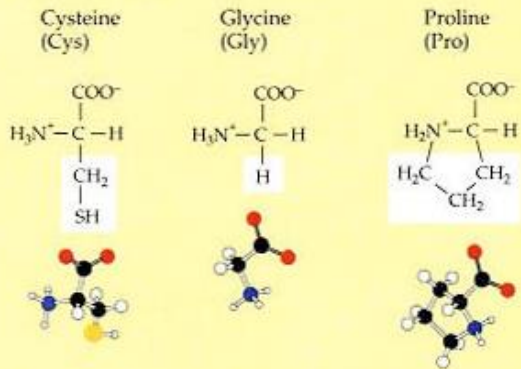
A. Amino acids with electrically charged side chains



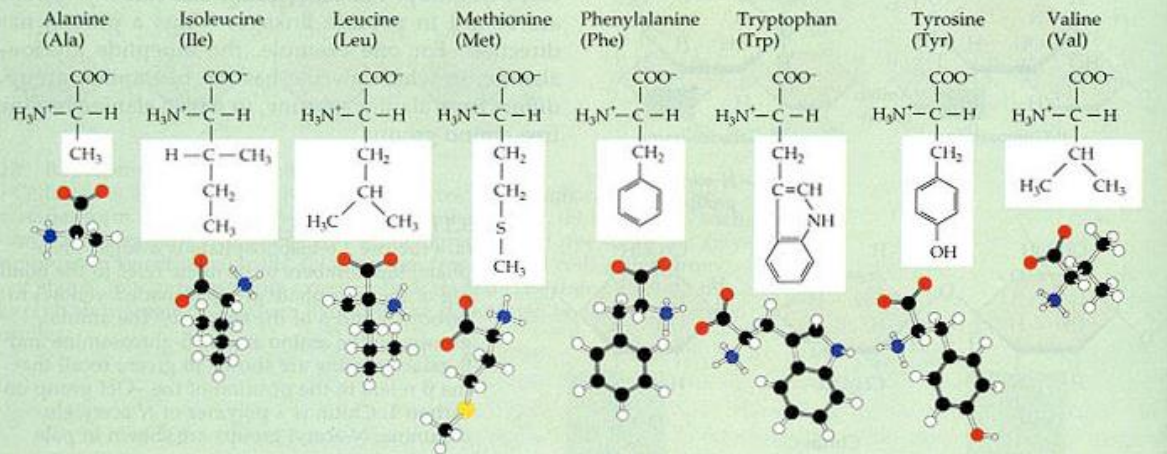
B. Amino acids with polar but uncharged side chains



C. Special cases



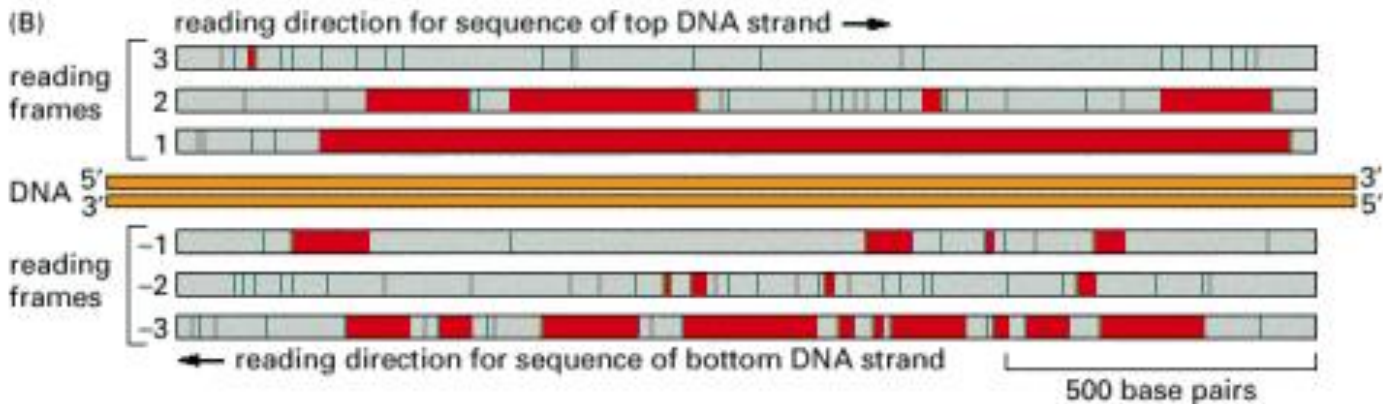
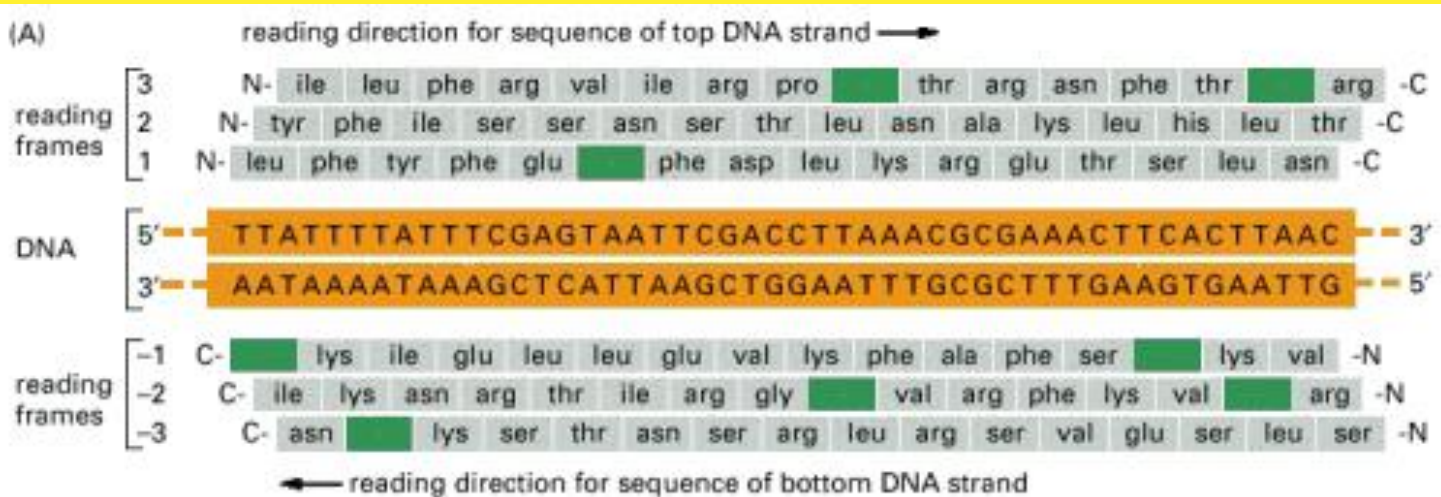
D. Amino acids with hydrophobic side chains



		Second letter				
		U	C	A	G	
First letter	U	UUU } Phe	UCU } Ser	UAU } Tyr	UGU } Cys	U C A G
		UUC } Phe	UCC } Ser	UAC } Tyr	UGC } Cys	
		UUA } Leu	UCA } Ser	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 } Ile	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	

Third letter

Any region of the DNA sequence can, in principle, code for six different amino acid sequences, because any one of three different reading frames can be used to interpret each of the two strands.



Protein folding

A human Hemoglobin:



How does it all look on a computer monitor?

A cDNA sequence

```
>gi|14456711|ref|NM_000558.3| Homo sapiens hemoglobin, alpha 1 (HBA1), mRNA  
ACTCTTCTGGTCCCCACAGACTCAGAGAGAACCCACCATGGTGCTGTCTCCTGCCGACAAGACCAA  
CGTCAAGGCCGCCTGGGGTAAGGTCGGCGCGCACGCTGGCGAGTATGGTGCGGAGGCCCTGGAG  
AGGATGTTCTGTCTTCCCCACCACCAAGACCTACTTCCCGCACTTCGACCTGAGCCACGGCTCT  
GCCCAGGTTAAGGGCCACGGCAAGAAGGTGGCCGACGCGCTGACCAACGCCGTGGCGCACGTGG  
ACGACATGCCCAACGCGCTGTCCGCCCTGAGCGACCTGCACGCGCACAAGCTTCGGGTGGACCCG  
GTCAACTTCAAGCTCCTAAGCCACTGCCTGCTGGTGACCCTGGCCGCCACCTCCCCGCCGAGTTC  
ACCCCTGCGGTGCACGCCTCCCTGGACAAGTTCCTGGCTTCTGTGAGCACCGTGCTGACCTCCAAA  
TACCGTTAAGCTGGAGCCTCGGTGGCCATGCTTCTTGCCCCTTGGGCCTCCCCCCAGCCCCTCCTC  
CCCTTCTGCACCCGTACCCCCGTGGTCTTTGAATAAAGTCTGAGTGGGCGGC
```


A cDNA sequence (reading frame)

>gi|14456711|ref|NM_000558.3| **Homo sapiens hemoglobin, alpha 1 (HBA1), mRNA**
ACTCTTCTGGTCCCCACAGACTCAGAGAGAACCCACC**ATG**GTGCTGTCTCCTGCCGACAAGACCA
ACGTCAAGGCCGCCTGGGGTAAGGTCGGCGCGCACGCTGGCGAGTATGGTGCGGAGGCCCTGGA
GAGGATGTTCTGTCTTCCCCACCACCAAGACCTACTTCCCGCACTTCGACCTGAGCCACGGCTC
TGCCAGGTTAAGGGCCACGGCAAGAAGGTGGCCGACGCGCTGACCAACGCCGTGGCGCACGTG
GACGACATGCCCAACGCGCTGTCCGCCCTGAGCGACCTGCACGCGCACAAGCTTCGGGTGGACCC
GGTCAACTTCAAGCTCCTAAGCCACTGCCTGCTGGTGACCCTGGCCGCCACCTCCCCGCCGAGTT
CACCCCTGCGGTGCACGCCTCCCTGGACAAGTTCCTGGCTTCTGTGAGCACCGTGCTGACCTCAA
ATACCGT**TAA**GCTGGAGCCTCGGTGGCCATGCTTCTTGCCCCTTGGGCCTCCCCCAGCCCCTCC
TCCCCTTCTGCACCCGTACCCCGTGGTCTTTGAATAAAGTCTGAGTGGGCGGC

A protein sequence

>gi|4504347|ref|NP_000549.1| **alpha 1 globin [Homo sapiens]**
MVLSPADKTNVKAAWGKVGAGHAGEYGAELERMFLSFPTTKTYFPHFDLSHGSAQVKGHGKKVA
DALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFL
ASVSTVLTSKYR

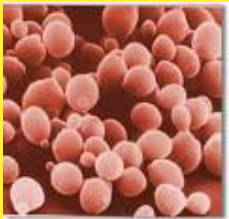
And, a whole genome...

ACTCTTCTGGTCCCCACAGACTCAGAGAGAACCCACCATGGTGCTGTCTCCTGCCGACAAGACCAACGTC
AAGGCCGCCTGGGGTAAGGTTCGGCGCGCACGCTGGCGAGTATGGTGCGGAGGCCCTGGAGAGGATGTT
CCTGTCCCTTCCCCACCACCAAGACCTACTTCCCAGCACTTCGACCTGAGCCACGGCTCTGCCAGGTTAAG
GGCCACGGCAAGAAGGTGGCCGACGCGCTGACCAACGCCGTGGCGCACGTGGACGACATGCCCAACGC
GCTGTCCGCCCTGAGCGACCTGCACGCGCACAAAGCTTCGGGTGGACCCGGTCAACTTCAAGCTCCTAAG
CCACTGCCTGCTGGTGACCCTGGCCGCCACCTCCCCGCCGAGTTCACCCCTGCGGTGCACGCCTCCCT
GGACAAGTTCCTGGCTTCTGTGAGCACCGTGCTGACCTCCAAATACCGTTAAGCTGGAGCCTCGGTGGCC
ATGCTTCTTGCCCCTTGGGCCTCCCCCAGCCCCTCCTCCCCTTCTGCACCCGTACCCCGTGGTCTTT
GAATAAAGTCTGAGTGGGCGGCACTCTTCTGGTCCCCACAGACTCAGAGAGAACCCACCATGGTGCTGTC
TCCTGCCGACAAGACCAACGTCAAGGCCGCCTGGGGTAAGGTTCGGCGCGCACGCTGGCGAGTATGGTGC
GGAGGCCCTGGAGAGGATGTTCTGTCTTCCCCACCACCAAGACCTACTTCCCGCACTTCGACCTGAGC
CACGGCTCTGCCAGGTTAAGGGCCACGGCAAGAAGGTGGCCGACGCGCTGACCAACGCCGTGGCGCA
CGTGGACGACATGCCCAACGCGCTGTCCGCCCTGAGCGACCTGCACGCGCACAAAGCTTCGGGTGGACCC
GGTCAACTTCAAGCTCCTAAGCCACTGCCTGCTGGTGACCCTGGCCGCCACCTCCCCGCCGAGTTCACC
CCTGCGGTGCACGCCTCCCTGGACAAGTTCCTGGCTTCTGTGAGCACCGTGCTGACCTCCAAATACCGTT
AAGCTGGAGCCTCGGTGGCCATGCTTCTTGCCCCTTGGGCCTCCCCCAGCCCCTCCTCCCCTTCTGCA
CCCGTACCCCGTGGTCTTTGAATAAAGTCTGAGTGGGCGGCACTCTTCTGGTCCCCACAGACTCAGAGA
GAACCCACCATGGTGCTGTCTCCTGCCGACAAGACCAACGTCAAGGCCGCCTGGGGTAAGGTTCGGCGCG
CACGCTGGCGAGTATGGTGCGGAGGCCCTGGAGAGGATGTTCTGTCTTCCCCACCACCAAGACCTACT
TCCCGCACTTCGACCTGAGCCACGGCTCTGCCAGGTTAAGGGCCACGGCAAGAAGGTGGCCGACGCGC
TGACCAACGCCGTGGCGCACGTGGACGACATGCCCAACGCGCTGTCCGCCCTGAGCGACCTGCACGCGC
ACAAGCTTCGGGTGGACCCGGTCAACTTCAAGCTCCTAAGCCACTGCCTGCTGGTGACCCTGGCCGCCCA
CCTCCCCGCCGAGTTCACCCCTGCGGTGCACGCCTCCCTGGACAAGTTCCTGGCTTCTGTGAGCACCGTG
CTGACCTCCAAATACCGTTAAGCTGGAGCCTCGGTGGCCATGCTTCTTGCCCCTTGGGCCTCCCCCAGC
CCCTCCTCCCCTTCTGCACCCGTACCCCGTGGTCTTTGAATAAAGTCTGAGTGGGCGGCGCCGTGGCG
CACGTGGACGACATGCCCAACGCGCTGTCCGCCCTGAGCGACCTGCACGCGCACAAAGCTTCGGGTGGAC
CCGGTCAACTTCAAGCTCCTAAGCCACTGCCTGCTGGTGACCCTGGCCGCCACCTCCCCGCCGAGTTC
CCCCTGCGGTGCACGCCTCCCTGGACAAGTTCCTGGCTTCTGTGAGCACCGTGCTGACCTCCAAATACCG
TTAAGCTGGAGCCTCGGTGGCCATGCTTCTTGCCCCTTGGGCCTCCCCCAGCCCCTCCTCCCCTTCTG

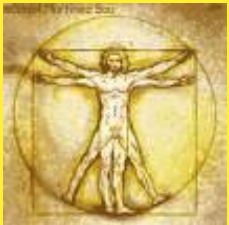
How big are whole genomes?



E. coli 4.6×10^6 nucleotides
– Approx. 4,000 genes



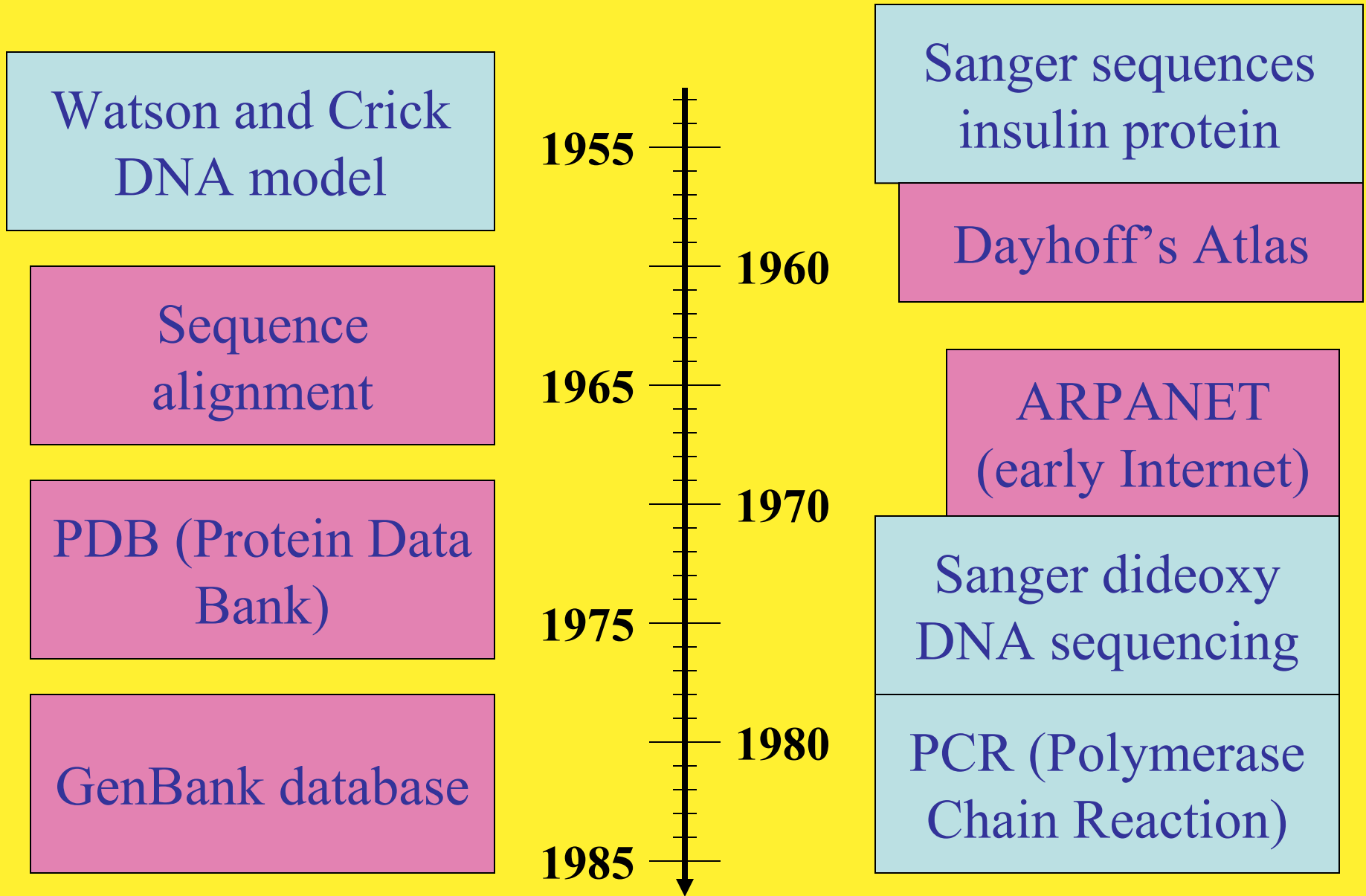
Yeast 15×10^6 nucleotides
– Approx. 6,000 genes

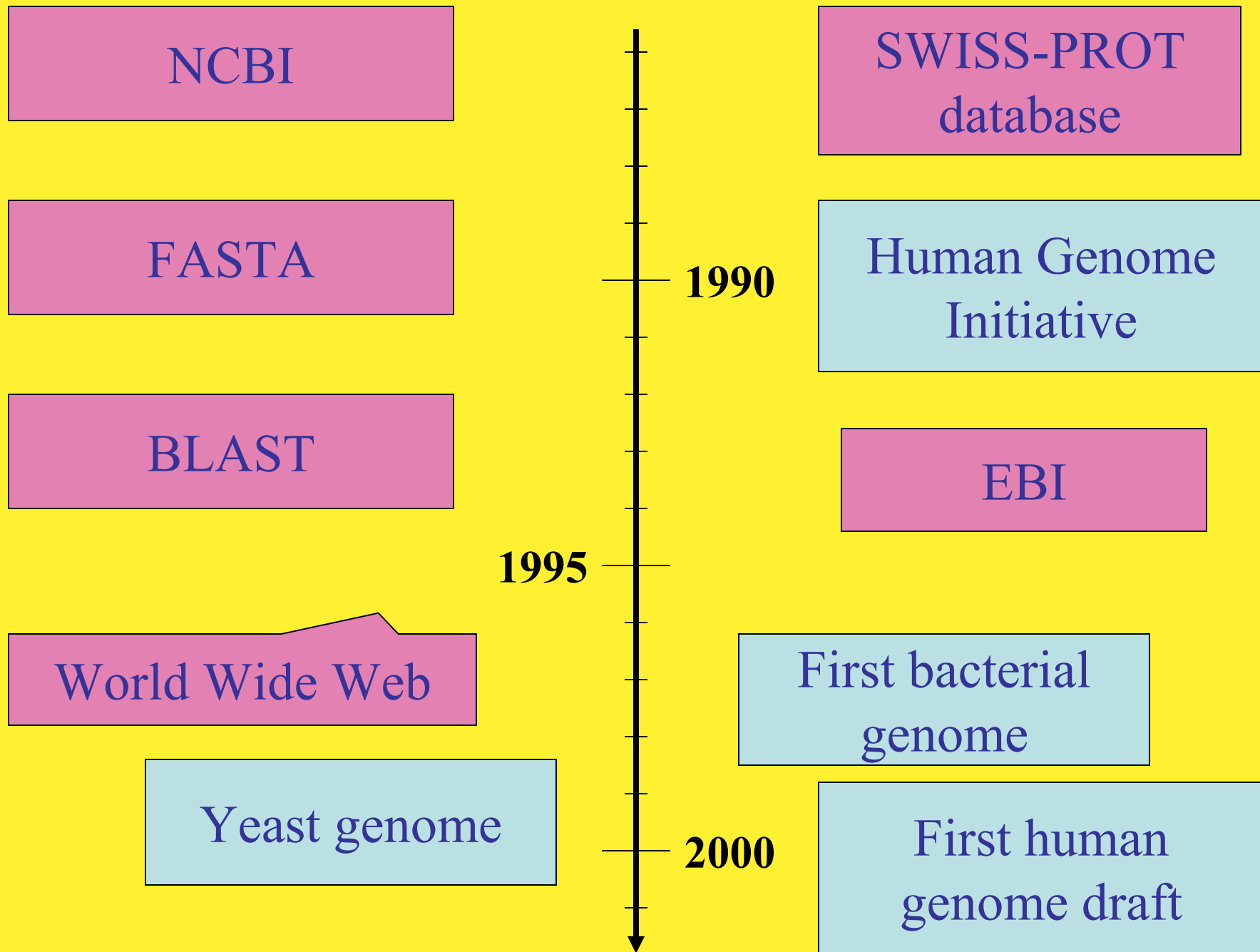


Human 3×10^9 nucleotides
– Approx. 30,000 genes

Smallest human chromosome 50×10^6 nucleotides

From DNA to Genome





Origin of bioinformatics and biological databases:

The first protein sequence reported was that of bovine insulin in **1956**, consisting of 51 residues.

Nearly a decade later, the first nucleic acid sequence was reported, that of yeast tRNA^{alanine} with 77 bases.

In 1965, **Dayhoff** gathered all the available sequence data to create the first bioinformatic database (*Atlas of Protein Sequence and Structure*).

The Protein DataBank followed in 1972 with a collection of ten X-ray crystallographic protein structures. The SWISSPROT protein sequence database began in 1987.

Complete Genomes

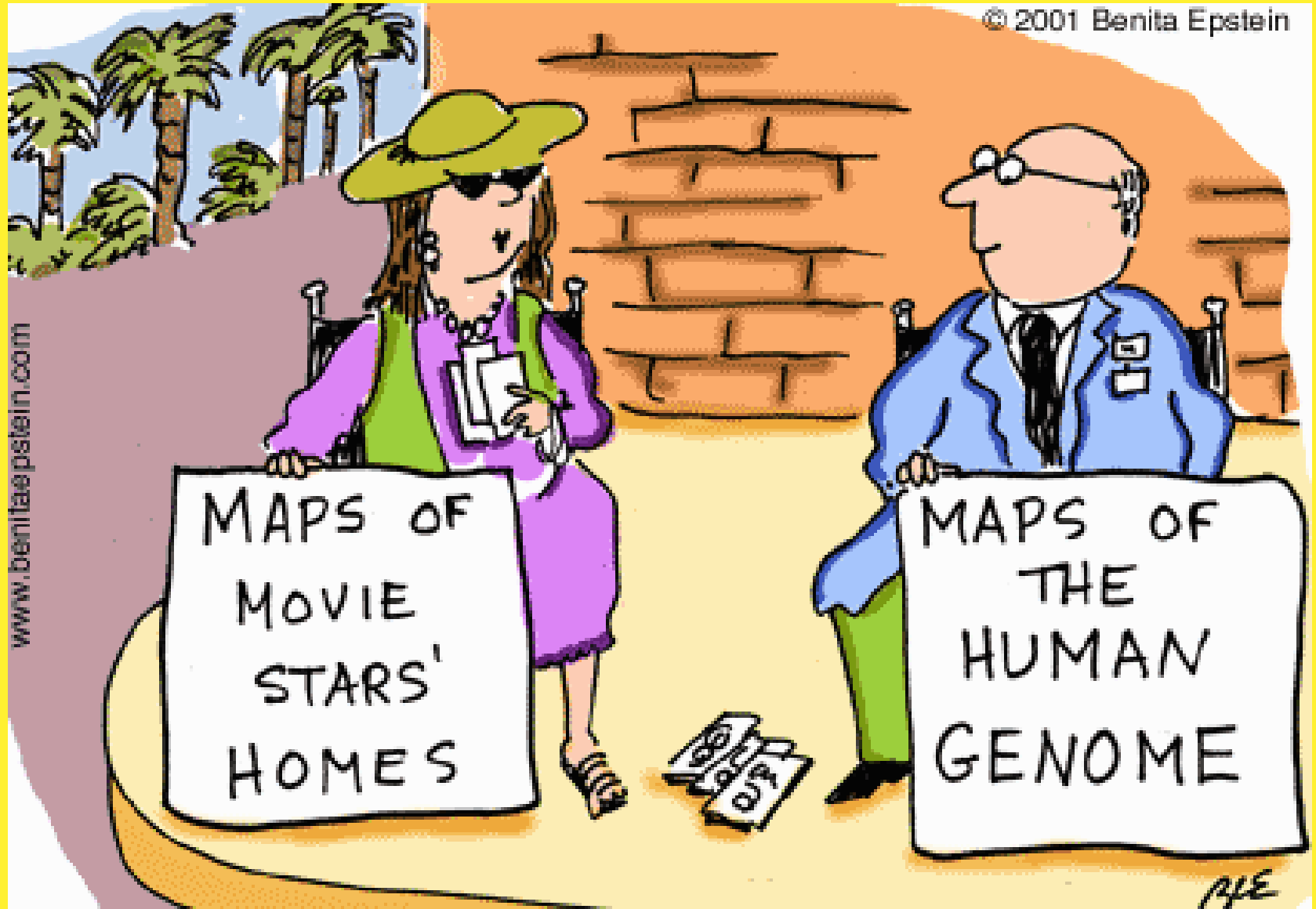
as of August 2011:

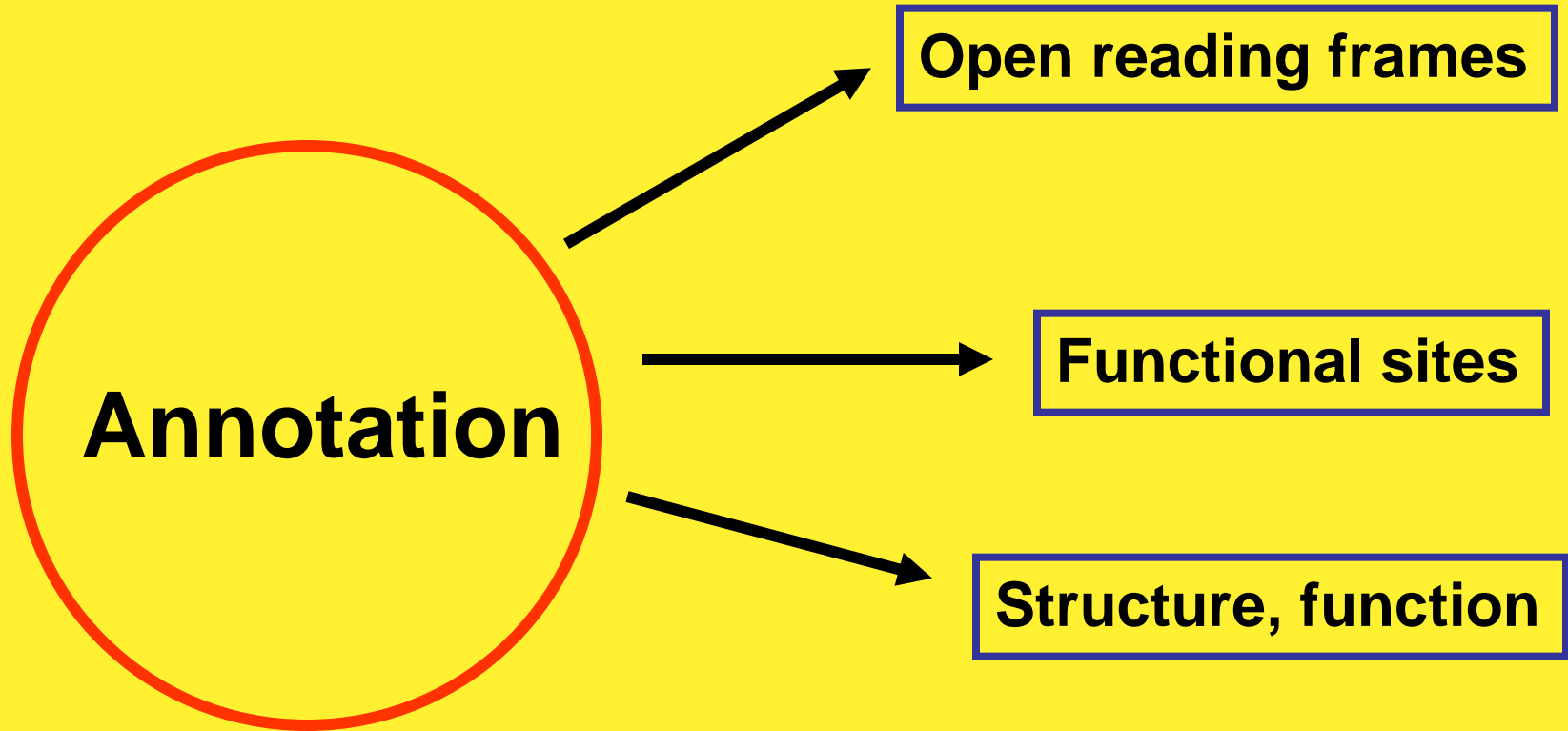
Eukaryotes 37

Prokaryotes 1708

Total 1745

What can we do with sequences and other type of molecular information?





**CCTGACAAATTCGACGTGCGGCATTGCATGCAGACGTGCATG
CGTGCAAATAATCAATGTGGACTTTTCTGCGATTATGGAAGAA
CTTTGTTACGCGTTTTTGTTCATGGCTTTGGTCCCGCTTTGTTC
AGAATGCTTTTAATAAGCGGGGTTACCGGTTTGGTTAGCGAGA
AGAGCCAGTAAAAGACGCAGTGACGGAGATGTCTGATG CAA
TAT GGA CAA TTG GTT TCT TCT CTG AAT**
..... TGAAAAACGTA

promoter

TF binding site

CCTGACAAATTCGACGTGC **GGCATTGCATGC** **AGACGTGCATG**
 CGTGCAA **TAATCA** ATGTGGACTTTTCTGC **GATTAT** GAAGA
 CTTTGTTACGCGTTTTTGTTCATGGCTTTGGTCCCGCTTTGTTTC
 AGAATGCTTTTAATAAGCGGGGTTACCGGTTTGGTTAGCGAGA
 AGAGCCAGTAAAAGACGCAGTGAC **GGAGAT** GTCTG **ATG CAA**
TAT GGA CAA TTG GTT TCT TCT CTG AAT
 TGAAAAACGTA

Transcription Start Site

Ribosome binding Site

ORF = Open Reading Frame
CDS = Coding Sequence

Comparative genomics

Comparing ORFs



Identifying orthologs



Inferences on structure and function

Comparing functional sites

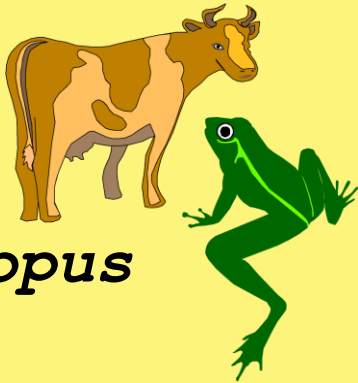


Inferences on regulatory networks

Alignment preproinsulin

Xenopus

Bos



MALWMQCLP-LVLVLLFSTPNTTEALANQHL

MALWTRLRPLLALLALWPPPPARAFVNQHL

**** : * * . * : * : . . * : . * : ****

Xenopus

Bos

CGSHLVEALYLVCGRGFFYYPKIKRDIEQ

CGSHLVEALYLVCGERGFFYTPKARREVEG

***** : ***** ** : * : : *

Xenopus

Bos

AQVNGPQDNELDG-MQFQPQEYQMKRGIV

PQVG---ALELAGGPGAGGLEGPPQKRGIV

. ** . ** * * *****

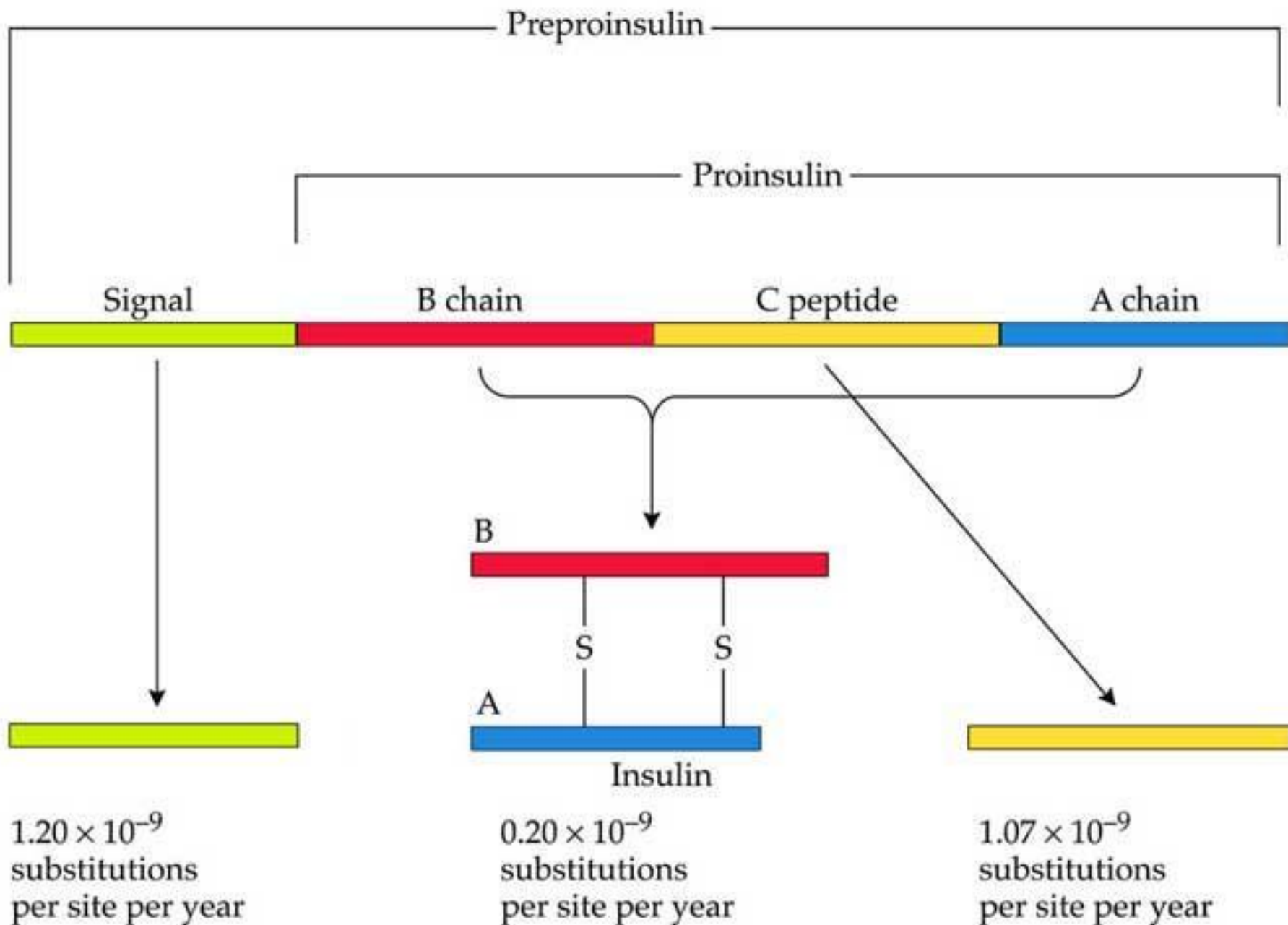
Xenopus

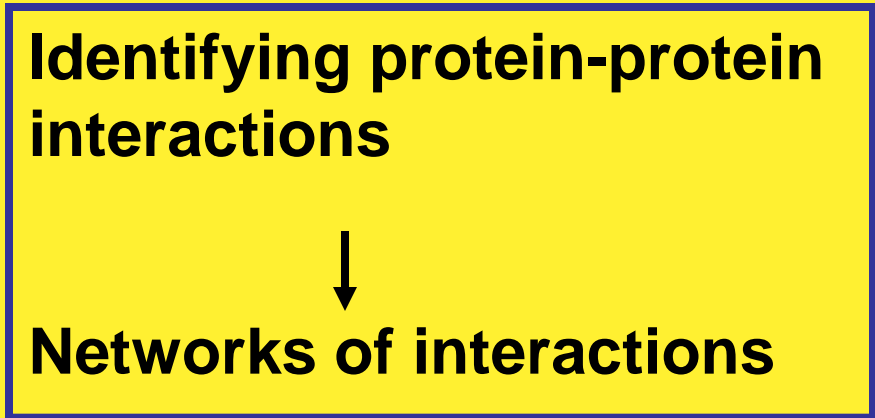
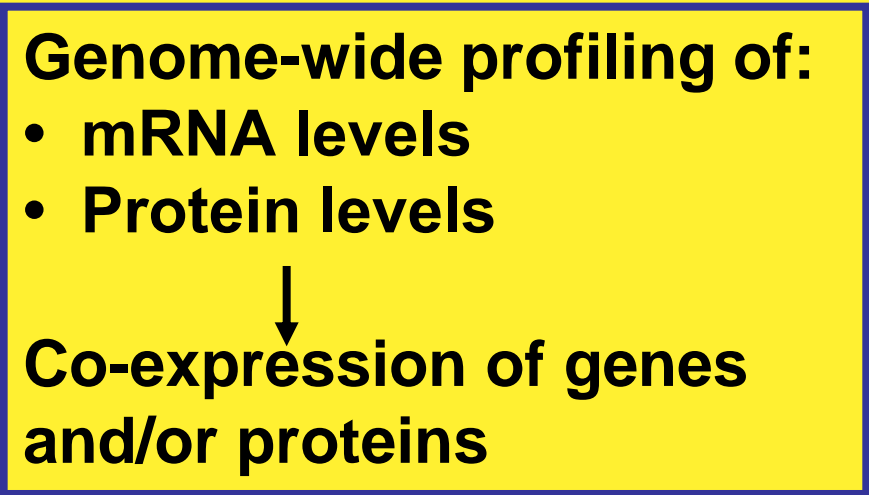
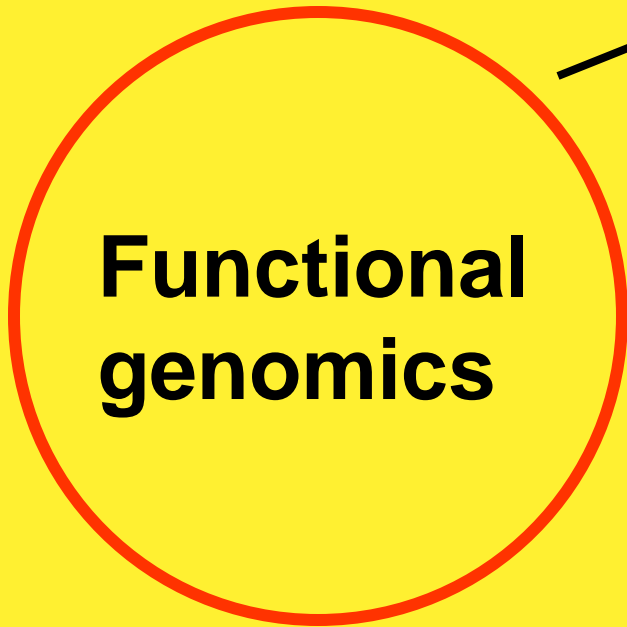
Bos

EQCCHSTCSLFQLENYCN

EQCCASVCSLYQLENYCN

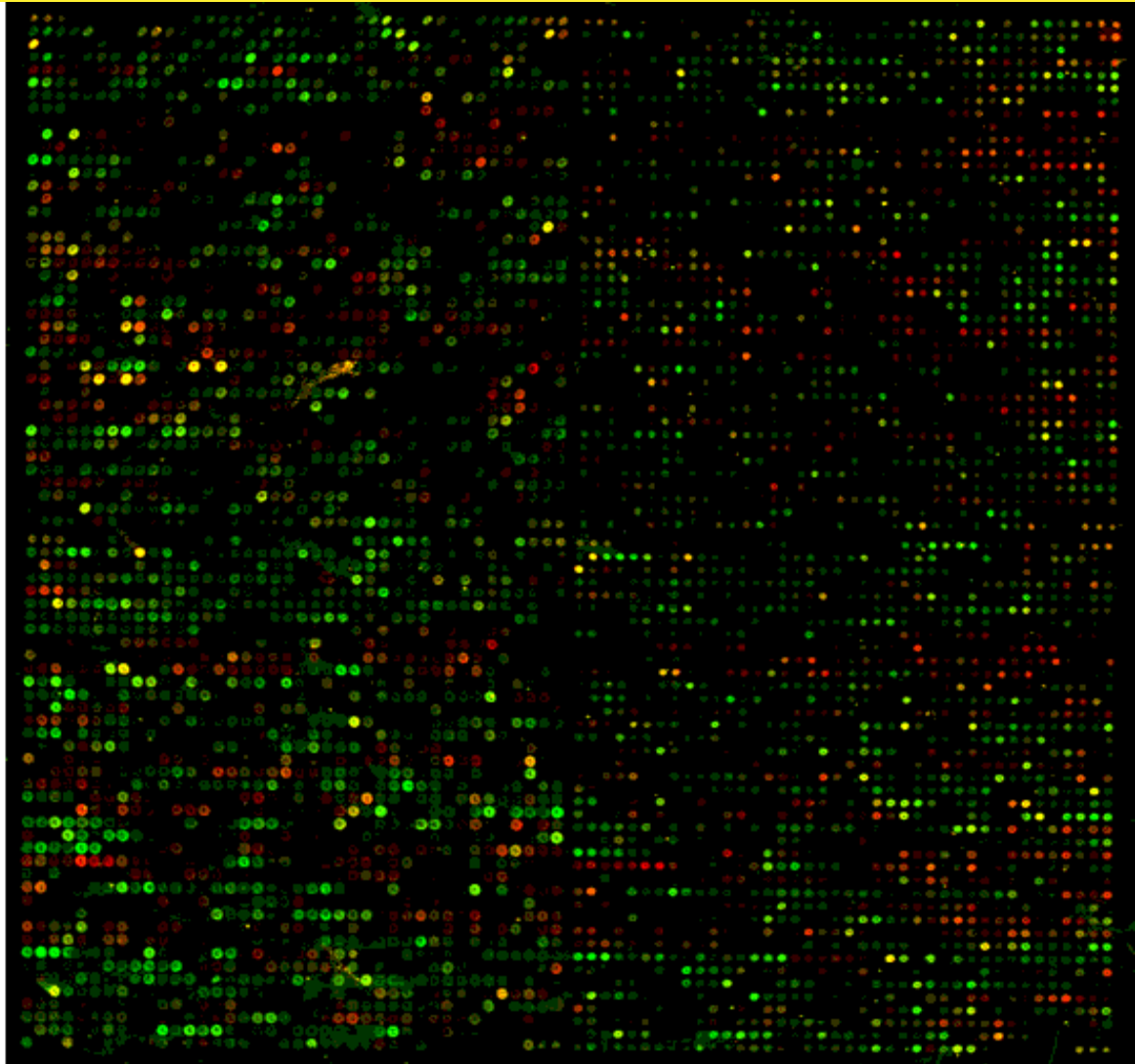
**** * . **** : *****





Understanding the function of genes and other parts of the genome

**The
complete
S. cerevisiae
genome
on a
microarray
chip
hybridised
to RNA
from cultures in
anaerobic and
aerobic
stationary phase**





**Structural
genomics**



**Assign structure to all
proteins encoded in
a genome**