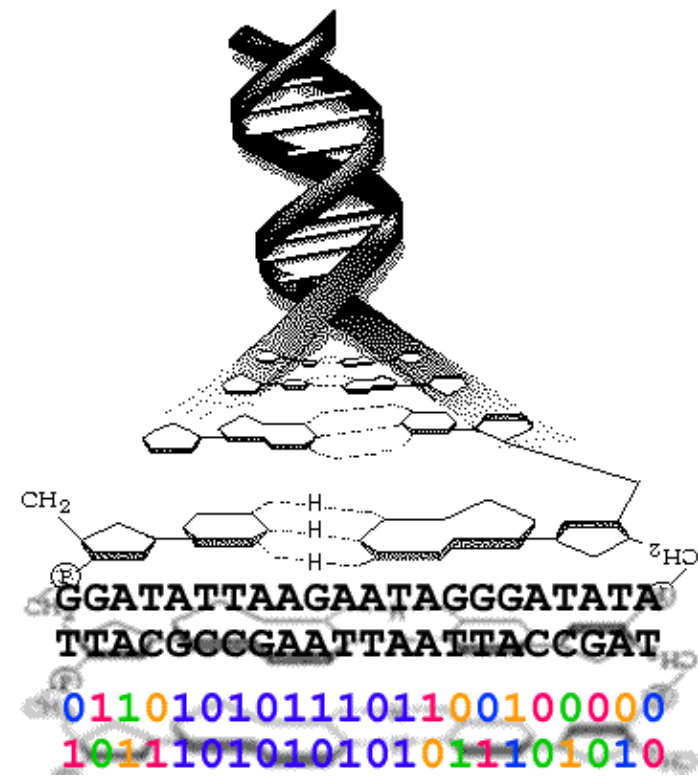
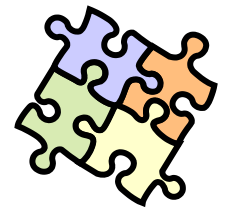


# *Biology & Big Data*

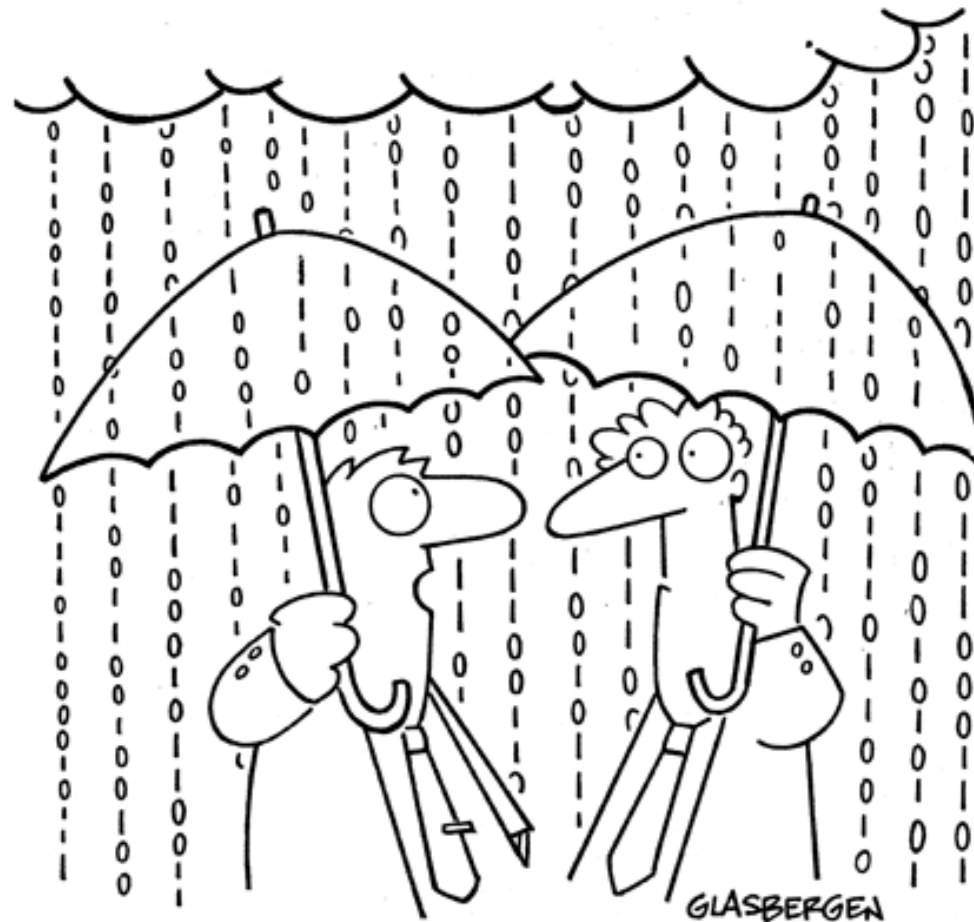


Debasis Mitra

Professor, Computer Science, FIT



# Cloud?



**"I don't know much about cloud computing,  
but I think it might be responsible for  
the strange weather we're having."**

---

## Data as Service

- **Transparent to user**
- **Multiple locations**
- **Robustness**
- **Cost-based speed**

## Software as Service

- **Software location transparent**
- **Pay-as-you-use**
- **Software running platform transparent**

***Underlying middle-layer makes everything transparent***

---

---

## Platform as Service

- **Middle-layer +**
- **Software environment +**
- **User interface**

## Infrastructure as Service

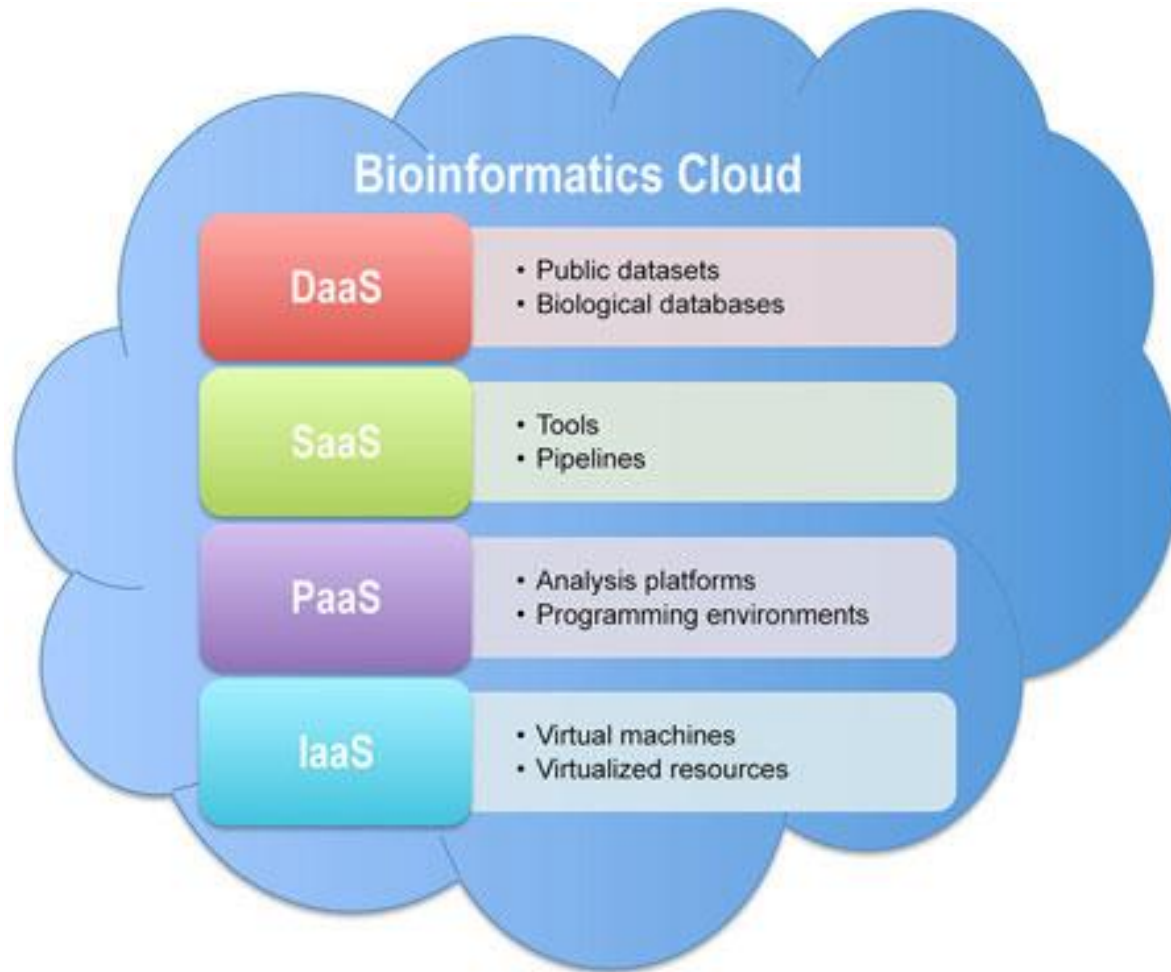
- **Hardware +**
- **Internet +**
- **Virtual Machine**

*Source: Dai et al. Biology Direct 2012, 7:43 <http://www.biology-direct.com/content/7/1/43>*

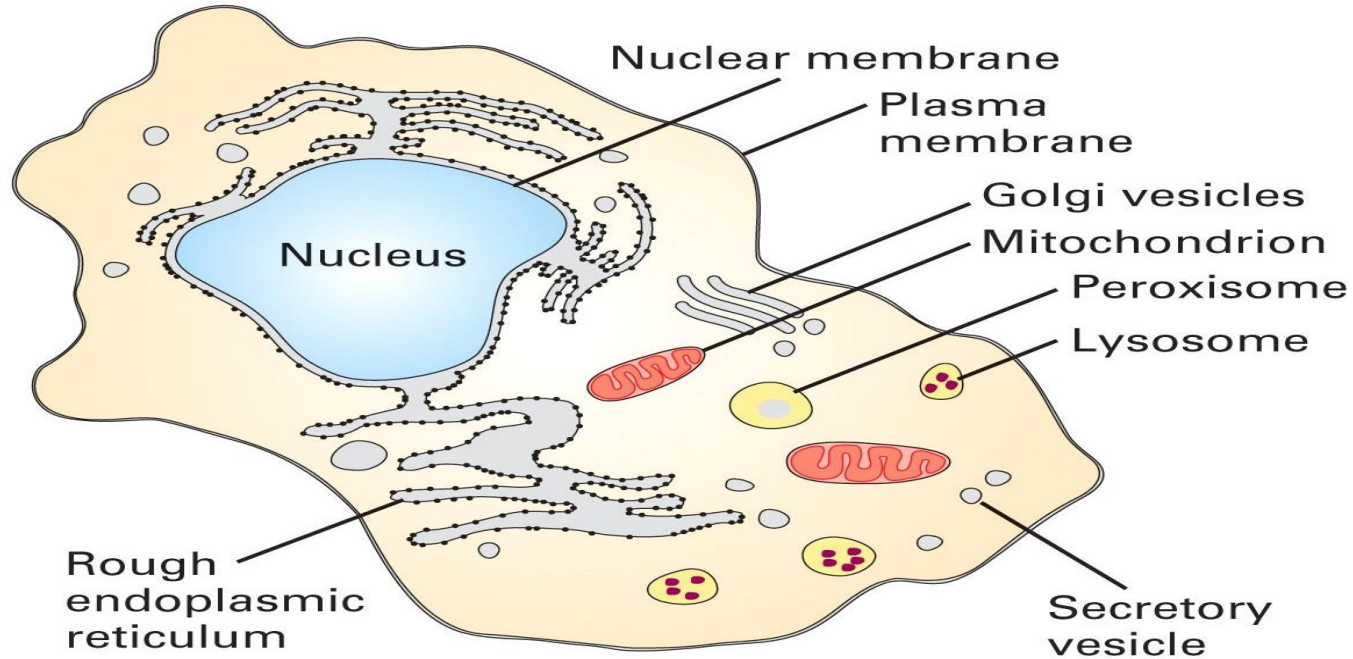
---

---

# *Bio-informatics Cloud / Grid*

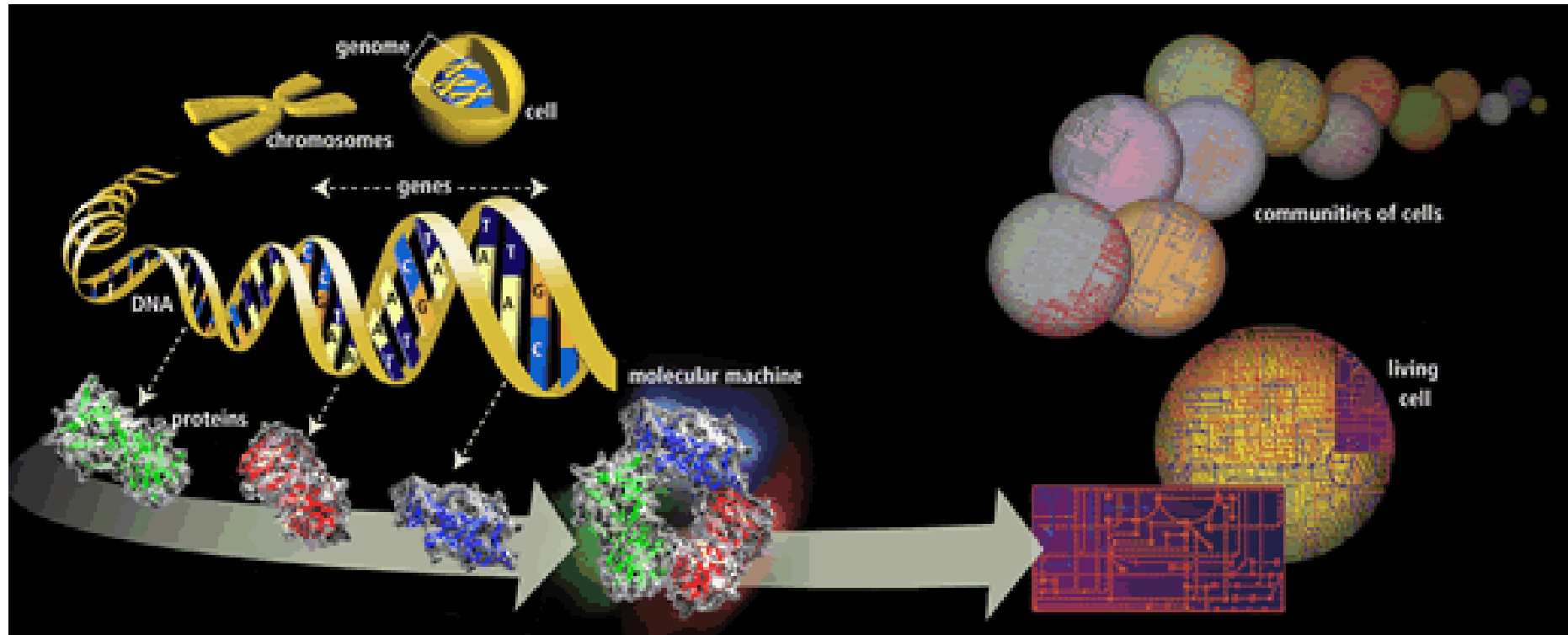


# Life begins with Cell



- A cell is a smallest structural unit of an organism that is capable of independent functioning
- All cells have some common features

# Molecular Biology



*An Introduction to Bioinformatics Algorithms (Computational Molecular Biology)*  
Neil C. Jones, Pavel A. Pevzner

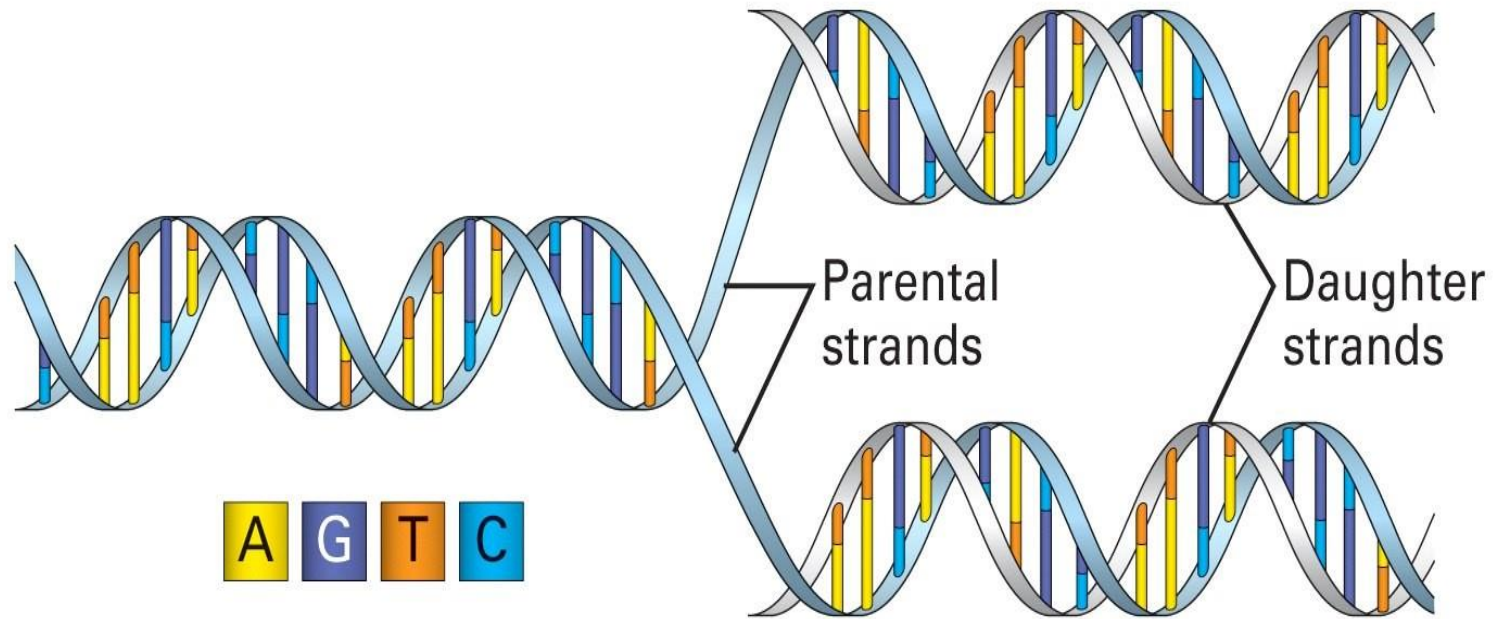
---

# Organizations of life

- **Nucleus = library**
  - **Chromosomes = bookshelves**
  - **Genes = books**
  - Same libraries and the same sets of books for every cell in an organism
  - Books represent all the information to carry out its various functions.
-

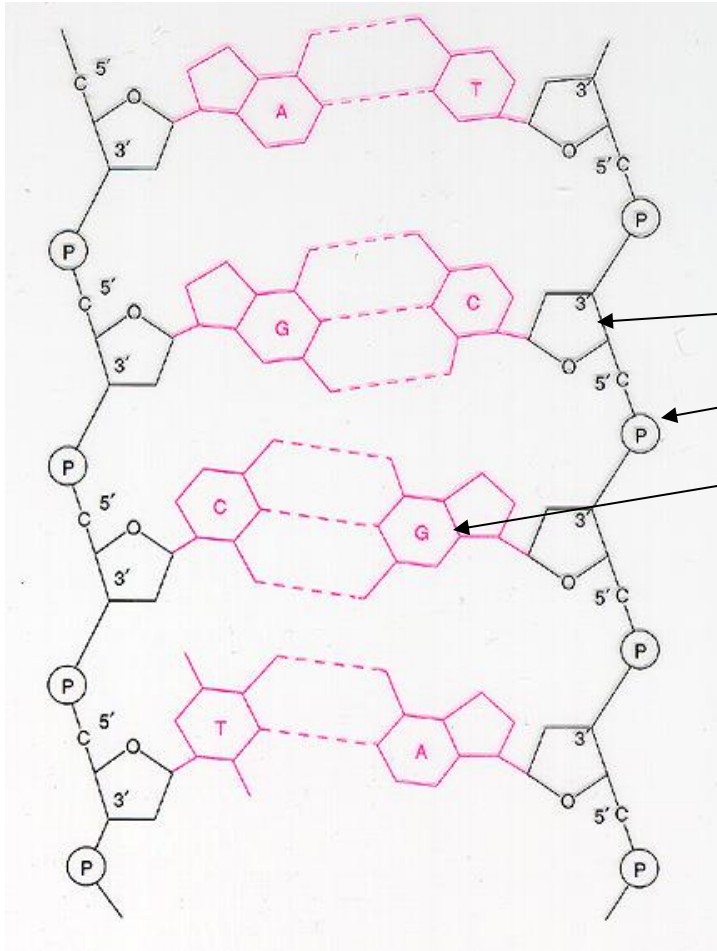


# DNA: The Code of Life



- The structure and the four genomic letters code for all living organisms
- Adenine, Guanine, Thymine, and Cytosine which pair A-T and C-G on complimentary strands.

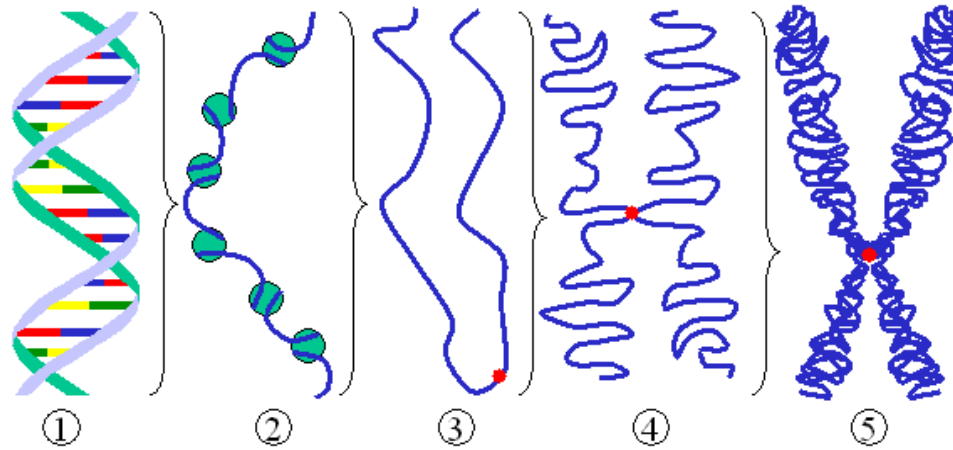
# DNA, continued



- DNA has a double helix structure which is composed of
  - sugar molecule
  - phosphate group
  - and a base (A,C,G,T)

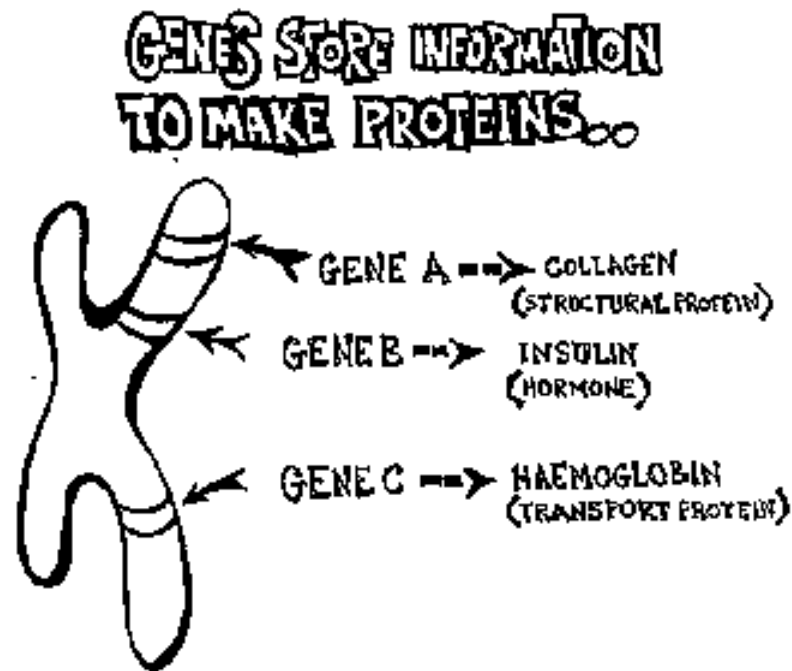
5' ATTTAGGCC 3'  
3' TAAATCCGG 5'

# Genetic Information: Chromosomes



# Genes Make Proteins

- genome -> genes -> protein (forms cellular structural & life functional) -> pathways & physiology



# Central Paradigm of Bioinformatics

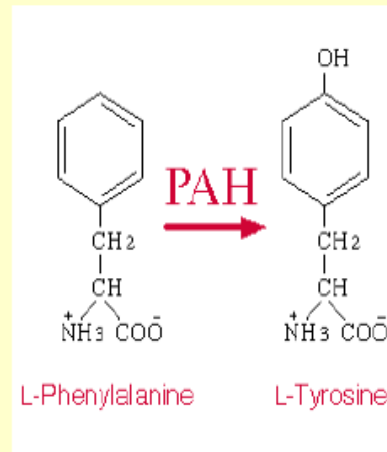
Genetic Information

MVHLTPEEKT  
AVNALWGKVN  
VDAVGGEALG  
RLLVVYPWTQ  
RFFESFGDLS  
SPDAVMGNPK  
VKAHGKKVLG  
AFSDGLAHL  
NLKGTFSQLS  
ELHCDKLHVD  
PENFRLLGNV  
LVCVLARNFG  
KEFTPQMCAA  
YQKVVAGVAN  
ALAHKYH

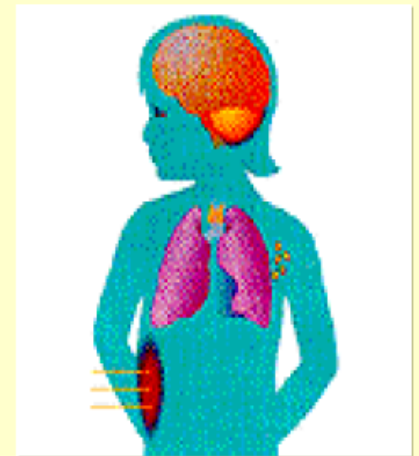
Molecular Structure



Biochemical Function



Phenotype (Symptoms)



---

## Some Terminology

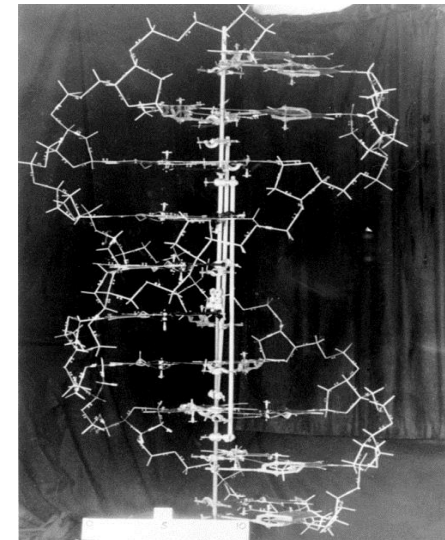
- **Genome**: an organism's genetic material
    - *Human genome is about 3,000,000,000 base-pair long*
  - **Gene**: a discrete units of hereditary information located on the chromosomes and consisting of DNA.
  - **Nucleic acid**: Biological molecules(RNA and DNA) that allow organisms to reproduce: A, T (U), C, G
-

# Major events in the history of Molecular Biology 1952 - 1960

- **1952-1953** James D. Watson and Francis H. C. Crick deduced the double helical structure of DNA with four molecules A, T, C, G
- 1 Biologist
- + 1 Physics student
- + 900 words
- -----
- = Nobel Prize



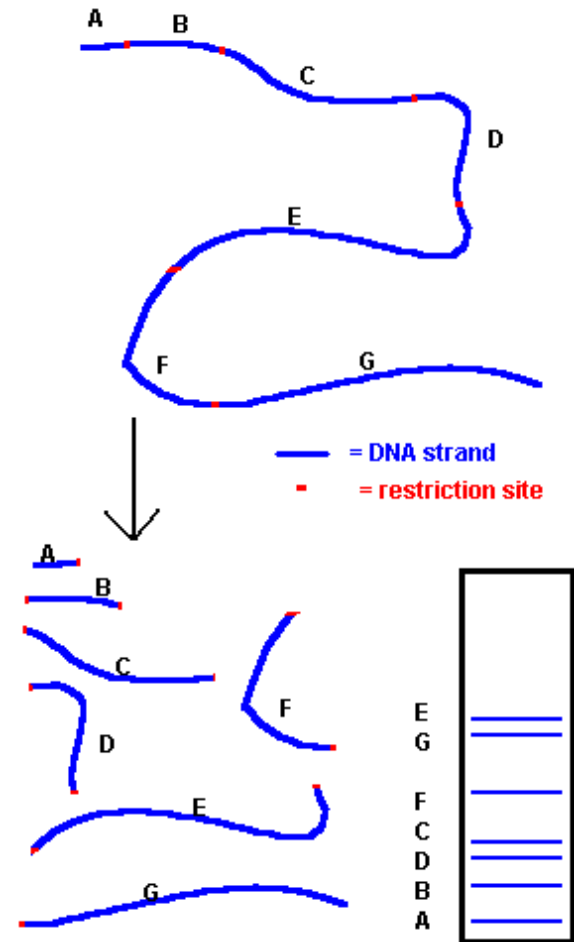
James Watson and Francis Crick



Original DNA model (scale: 1 cm = 100 Angstroms)  
Cold Spring Harbor Laboratory Archives

# Major events in the history of Molecular Biology 1970

- 1970 Howard Temin and David Baltimore independently isolate the first restriction enzyme
- DNA can be cut into reproducible pieces with restriction enzymes;  
(gene cloning or recombinant DNA technology)





# Major events in the history of Molecular Biology

## 1986 - 1995

- **1986** Leroy Hood: Developed automated sequencing mechanism
- **1986** Human Genome Initiative announced
- **1990** The 15 year Human Genome project is launched by congress
- **1995** Moderate-resolution maps of chromosomes 3, 11, 12, and 22 maps published (These maps provide the locations of “markers” on each chromosome to make locating genes easier)



Leroy Hood



---

# Major events in the history of Molecular Biology

## 1995-1996

- **1995** John Craig Venter: First **bacterial genomes** sequenced
- Challenged the genome sequencing project by developing ‘*shotgun*’ approach –
- approach depends on assembling the sequences by **computer**



John Craig Venter

---

---

# Major events in the history of Molecular Biology

## 1997 - 1999

- **1997** E. Coli sequenced
  - **1998** PerkinsElmer, Inc.. Developed 96-capillary sequencer
  - **1998** Complete sequence of the **Caenorhabditis elegans genome**
  - **1999** First human chromosome (number 22) sequenced
-

# Major events in the history of Molecular Biology 2000-2001

- **2000** Complete sequence of the euchromatic portion of the *Drosophila melanogaster* genome
- **2001** International Human Genome Sequencing: first draft of the sequence of the human genome published



---

# Major events in the history of Molecular Biology

## 2003- Present

- **April 2003** Human Genome Project Completed. Mouse genome is sequenced.
- **Jan 15, 2014, *Illumina*** Your genome may be sequenced < \$1,000
- **IBM Challenge:** <\$100



---

*Why sequence?*  
MUTAsHONS

- What happens to genes when the DNA is mutated?

Normal DNA sequence:

**ATCTAG**

Mutated DNA sequence:

**ATC**G**AG**



---

# The Good, the Bad, and the Silent

- Mutations can serve the organism in three ways:

A mutation can cause a trait that enhances the organism's function:

- **The Good :** Mutation in the sickle cell gene provides resistance to malaria.

A mutation can cause a trait that is harmful, sometimes fatal to the organism:

- **The Bad :** Huntington's disease, a symptom of a gene mutation, is a degenerative disease of the nervous system.

- **The Silent:** A mutation can simply cause no difference in the function of the organism.

# Real data: human & fruitfly eyeless

- This is a global alignment of human & fruitfly Eyeless-gene

Next few slides are from Dr Avril Coghlan, Sanger Institute, UK

```
human/1-422 .....
fly/1-898 1 MFTLQPTPTAIGTVVPPWSAGTLIERLPSLEDMAHKDNVIAMRNLPCLGT 50

human/1-422 1 .....MQNSHSG 7
fly/1-898 51 AGGSGGLGGIAGKPSPTMEAVEASTASHPHSTSSYFATYYHLTDDECHSG 100

human/1-422 8 VNQLGGVFNGRPLPDSTRQKIVELAHSGARPCDISRILQVSNCGVSKIL 57
fly/1-898 101 VNQLGGVFNVGRPLPDSTRQKIVELAHSGARPCDISRILQVSNCGVSKIL 150

human/1-422 58 GRYYETGSIRPRAIGGSKPRVATPEVVSKIADYKRECPSIFAWEIRDRL 107
fly/1-898 151 GRYYETGSIRPRAIGGSKPRVATPEVVSKISQYKRECPSIFAWEIRDRL 200

human/1-422 108 SEGVCTNDNIPSVSSINRVLRLNLA SEKQQM ..... 137
fly/1-898 201 QENVCTNDNIPSVSSINRVLRLNLA AQKEDQSTGSGSSSTSA GNSISAKVS 250

human/1-422 138 .....GA ..... DG 141
fly/1-898 251 VSI GGNVSNVAVSGSRGTLSSSTDLMQTATPLN SSESGGASNSGEGSEQEA 300

human/1-422 142 MYDKLRMLNGQTG ..... 154
fly/1-898 301 IYEKLRLLNTQHAAGPGLPARAAPLVGQSPNHLGTRSSHPQLVHGHNQ 350

human/1-422 155 ..... SWGTR...PGWYFGTSVPGQPTQ ..... 174
fly/1-898 351 ALQQHQQQSWPPRHYSGSWYP-TSLSEIPISSAPNIASVTAYASGPSLAH 399

human/1-422 175 ..... DG CQQQE...GGGENTM 188
fly/1-898 400 SLSPNDIESLASIGHQRNCPVATEDIHLKKELDGHQSD ETGSGEGENSN 449

human/1-422 189 SISSNGEDSDEAQMRLQLKRRKLRNRTSFTQEQIEALEKEFERTHYPDV 238
fly/1-898 450 GGASNI GNTEDDQARLILKRRKLRNRTSFTNDQIDSLEKEFERTHYPDV 499

human/1-422 239 ARERLAAKIDLPEARIQVWFSNRRAKWRREEKLRNQRRQASNTPSHIPIS 288
fly/1-898 500 ARERLAGKIGLPEARIQVWFSNRRAKWRREEKLRNQRRTPNSTGASATSS 549

human/1-422 289 SFSSTSVYQIPQPPTTVPSSFTSGSMLGRDALTNTYSALPPMPSTMA 338
fly/1-898 550 STSATASLTDSPNSLSACSLLSGSAGGSPSVSTINGLSS ..... PSTLST 594

human/1-422 339 M-NLP .....MQPPVPSQTSSYSCMLPTSPSVNGRSYD ..... TYT 373
fly/1-898 595 NVNAPT LGAGIDSSSESTPIPHIRPSC---TSDNDNGRQSEDCRRVCSPC 641

human/1-422 374 PPHMQTHMNSQPMGTSSTTSTGLISFGVSVPVQVPGSEPDMSQYWPRLQ 422
fly/1-898 642 PLGVGGHQNTHHIQSNGHAQGGHALVPAIS ..... PRLNF 675

human/1-422 .....
fly/1-898 676 NSGSFGAMYSNMHHTALSMSDSYGAVTIPSFNHSAVGPLAPPSPIPQQG 725

human/1-422 .....
fly/1-898 726 DLTPSSLYPCHMTLRPPMAPAHHHIVPGDGGRPAGVGLGSGQSANLGAS 775

human/1-422 .....
fly/1-898 776 CSGSGYEVL SAYALPPPMASSSAADSSSFAASSASANVTPHHTIAQESC 825

human/1-422 .....
fly/1-898 826 PSPCSSASHFGVAHSSGFSDDPISPAVSSYAHMSYNYASSANTMTPTSSAS 875

human/1-422 .....
fly/1-898 876 GTSAHVAPGKQQFFASCFYSPWW 898
```



# Real data: human & fruitfly eyeless

```
human/1-422 .....
fly/1-898 1 MFTLQPTPTAIGTVVPPWSAGTLIERLPSLEDMAHKDNV I AMRNLPCLGT 50

human/1-422 1 .....MQNS HSG
fly/1-898 51 AGGSGGLGG I AGKPSPTMEAVEASTASHPHSTSSYFATYYHLTDDECHSG

human/1-422 8 VNQLGGVFVNGRPLPDSTRQKIVELAHSGARPCD I SRILQVSNCGVSK I L
fly/1-898 101 VNQLGGVFVGGRPLPDSTRQKIVELAHSGARPCD I SRILQVSNCGVSK I L

human/1-422 58 GRYYETGS I RPRAIGGSKPRVATPEVVSKI AQYKRECPSIFAW E I DRLL
fly/1-898 151 GRYYETGS I RPRAIGGSKPRVATAEVVSKI SQYKRECPSIFAW E I DRLL

human/1-422 108 SEGVCTNDN I PSVSS I NRVLRLNLA SEK QQM .....
fly/1-898 201 QENVCTNDN I PSVSS I NRVLRLNLA AQK EDQSTGSGSSSTSAGNS I SAKVS

human/1-422 138 .....GA ..... DG 141
fly/1-898 251 VSIGGNVSNVAVSGSRGTLSSSTDLMQTATPLNSSESGGA S NSGEGSEQE A 300

human/1-422 142 MYDKLRMLNGQTG .....
fly/1-898 301 IYEKLRLLNTQHAAGPGLPARAAPLVGQSPNHLGTRSSHPQLVHGNNHQ 350

human/1-422 155 ..... SWGTR... PGWYFGTSPVGPQTQ .....
fly/1-898 351 ALQHQHQQQSWPPRHYSGSGWYP TSLSEIP I ISSAPN IASVTAYASG PSLAH 399

human/1-422 175 ..... DG CQQQE ..... GG GENTM 188
fly/1-898 400 SLSPND I ESLAS IGHQRNCPVATED IHLKKELDGHQSD E TGS GEGEN S 449

human/1-422 189 S I S S N G E D S D E A Q M R L Q L K R K L Q R N R T S F T Q E I E A L E K E F E R T H Y P D V F
fly/1-898 450 G G A S N I G N T E D D Q A R L I L K R K L Q R N R T S F T N D Q I D S L E K E F E R T H Y P D V F

human/1-422 239 A R E R L A A K I D L P E A R I Q V W F S N R R A K W R R E E K L R N Q R R Q A S N T P S H I P I S
fly/1-898 500 A R E R L A G K I G L P E A R I Q V W F S N R R A K W R R E E K L R N Q R R T P N S T G A S A T S S

human/1-422 289 S S F S T S V Y Q I P Q P T T P V S S F T S G S M L G R T D T A L T N T Y S A L P P M P S F T M A
fly/1-898 550 S T S A T A S L T D S P N S L S A C S S L L S G S A G G P S V S T I N G L S S ..... P S T L S T 594

human/1-422 339 M N L P ..... M Q P P V P S Q T S S Y S C M L P T S P S V N G R S Y D ..... T Y T 373
fly/1-898 595 M V N A P T L G A G I D S S E S E P T I P H I R P S C ... T S D N D N G R Q S E D C R R V C S P C 641

human/1-422 374 P P H M Q T H M N S Q P M G T S G T T S T G L I S P G V S V P V Q V P G S E P D M S Q Y W P R L Q
fly/1-898 642 P L G V G G H Q N T H H I Q S N G H A Q G H A L V P A I S ..... P R L N F 675

human/1-422 .....
fly/1-898 676 N S G S F G A M Y S N M H H T A L S M S D S Y G A V T P I P S F N H S A V G P L A P P S P I P Q Q G 725

human/1-422 .....
fly/1-898 726 D L T P S S L Y P C H M T L R P P P M A P A H H I V P G D G G R P A G V G L G S G Q S A N L G A S 775

human/1-422 .....
fly/1-898 776 C S G S G Y E V L S A Y A L P P P P M A S S S A A D S S F S A A S S A S A N V T P H H T I A Q E S C 825

human/1-422 .....
fly/1-898 826 P S P C S S A S H F G V A H S S G F S S D P I S P A V S S Y A H M S Y N Y A S S A N T M T P S S A S 875

human/1-422 .....
fly/1-898 876 G T S A H V A P G K Q Q F F A S C F Y S P W W 898
```

There are 2 short regions of high similarity

Outside those regions, there are many mismatches and gaps

It might be more sensible to make local alignments of one or both of the regions of high similarity

# Real data: human & fruitfly

```
human/1-398 1 HSGVNLGGVFN GRPLPDSTRQKIVELAHSGARPCDISRILQVSNCGVSV 50
fly/1-573 1 HSGVNLGGVFG GRPLPDSTRQKIVELAHSGARPCDISRILQVSNCGVSV 50

human/1-398 51 KILGRYYETGSIRPRAIGGSKPRVATPEVVSKI AQYKRECPSIF AWEIRD 100
fly/1-573 51 KILGRYYETGSIRPRAIGGSKPRVATAEVVSKISQYKRECPSIF AWEIRD 100

human/1-398 101 RLLSEG VCTNDNIPSVSSINRVLRLNLA SEKQQM..... 133
fly/1-573 101 RLLQEN VCTNDNIPSVSSINRVLRLNLA AQKEQSTGSGSSSTSAGNSISA 150

human/1-398 134 ..... GA..... 135
fly/1-573 151 KVSVSIGGNVSNVASGSRGTLSSSTDLMQTATPLNSSESGGA SNSGEGSE 200

human/1-398 138 - DGM YDKLRMLNGQTG..... 150
fly/1-573 201 QEA IYEKLRLLNTQHAAGPGPLEPARAAPLVGQSPNHLGTRSSHPQLVHG 250

human/1-398 151 ..... SWGTR... PGWYPGTSVPGQPTQ..... 170
fly/1-573 251 NHQALQQHQQQ SWPPRHYSGSWYP-TSLSEIPISSAPNIASVTAYASGPS 299

human/1-398 171 ..... DG CQQQE... GGE 181
fly/1-573 300 LAHSLSPNDIESLASIGHQRNCPVATEDIHLKKELDGHQSD ETGSGEGE 349

human/1-398 182 N T N S I S S N G E D S D E A Q M R L Q L K R K L Q R N R T S F T Q E Q I E A L E K E F E R T H Y P 231
fly/1-573 350 N S N G G A S N I G N T E D D Q A R L I L K R K L Q R N R T S F T N D Q I D S L E K E F E R T H Y P 399

human/1-398 232 D V F A R E R L A A K I D L P E A R I Q V W F S N R R A K W R R E E K L R N Q R R Q A S N T P S H I 281
fly/1-573 400 D V F A R E R L A G K I G L P E A R I Q V W F S N R R A K W R R E E K L R N Q R R T P N S T G A S A 449

human/1-398 282 P I S S S F S T S V Y Q P I P Q P T T P V S S F T S G S M L G R T D T A L T N T Y S A L P P M P S F 331
fly/1-573 450 T S S S T S A T A S L T D S P N S L S A C S S L L S G S A G G P S V S T I N G L S S . . . . . P S T 494

human/1-398 332 T M A N . N L P . . . . . M Q P P V P S Q T S S Y S C M L P T S P S V N G R S Y D . . . . . 366
fly/1-573 495 L S T N V N A P T L G A G I D S S E S P T P I P H I R P S C . . . T S D N D N G R Q S E D C R R V C 541

human/1-398 367 T Y T P P H M Q T H M N S Q P M G T S G T T S T G L I S P G V S 398
fly/1-573 542 S P C P L G V G G H Q N T H H I Q S N G H A Q G H A L V P A I S 573
```

This is a local alignment of human & fruitfly

What parts of the sequences were used in the local alignment?

# The Smith-Waterman algorithm

- **local alignment** of 2 sequences

The **alignment of all possible subsequences (parts) of sequences  $S_1$  and  $S_2$**

**The 0<sup>th</sup> row and 0<sup>th</sup> column of  $T$**  are first filled with zeroes

The recurrence relation used to fill table  $T$  is:

$$T(i, j) = \max \begin{cases} T(i-1, j-1) + \sigma(S_1(i), S_2(j)) \\ T(i-1, j) + \text{gap penalty} \\ T(i, j-1) + \text{gap penalty} \\ 0 \end{cases}$$

The traceback starts at the **highest scoring cell** in the matrix  $T$ , and travels up/left **while the score is still positive**

---





You fill in the whole of  $T$ , recording the previous cell (if any) used to calculate the value of each  $T(i, j)$ :

		G	G	C	T	C	A	A	T	C	A
	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0	0	0	0	2	2	0	0	2
C	0	0	0	2	0	2	0	1	1	2	0
C	0	0	0	2	1	2	1	0	0	3	1
T	0	0	0	0	4	2	1	0	2	1	2
A	0	0	0	0	2	3	4	3	1	1	3
A	0	0	0	0	0	1	5	6	4	2	3
G	0	2	2	0	0	0	3	4	5	3	1
G	0	2	4	2	0	0	1	2	3	4	2

The traceback starts at the **highest scoring cell** in the matrix  $T$ , and travels up/left **while the score is still positive**

		G	G	C	T	C	A	A	T	C	A
	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0	0	0	0	2	2	0	0	2
C	0	0	0	2	0	2	0	1	1	2	0
C	0	0	0	2	1	2	1	0	0	3	1
T	0	0	0	0	4	2	1	0	2	1	2
A	0	0	0	0	2	3	4	3	1	1	3
A	0	0	0	0	0	1	5	6	4	2	3
G	0	2	2	0	0	0	3	4	5	3	1
G	0	2	4	2	0	0	1	2	3	4	2

Work out the best local alignment from the traceback:

```

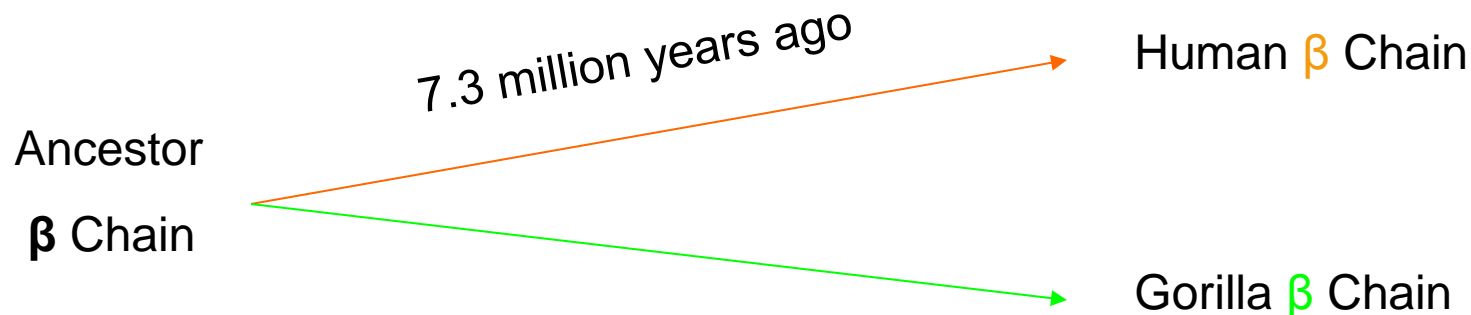
C T C A A
| |   | |
C T - A A

```

Score of the alignment is in the **bottom right cell of the traceback** ( $6 = 4 \times (\text{score of 2 per match}) + 1 \times (-2 \text{ per gap})$ )

# Molecular Clock

- Linus and Pauling found that  $\alpha$ -chains of human and gorilla differ by 2 residues, and  $\beta$ -chains by 1 residue.
- They then calculated the time of divergence between human and gorilla using evolutionary molecular clock.
- Gorilla and human  $\beta$  chain were found to diverge about 7.3 million years ago.





# Beta globins:

- Beta globin chains of closely related species are highly similar:
- Observe simple alignments below:

Human  $\beta$  chain: MVHLTPEEKSAVTALWGKV NVDEVGGEALGRLL

Mouse  $\beta$  chain: MVHLTDAEKAAVNGLWGKVNPDVVGGEALGRLL

Human  $\beta$  chain: VVYPWTQRFVESFGDLSTPDVVMGNPKVKAHGKKV LG

Mouse  $\beta$  chain: VVYPWTQRYFDSFGDLS SASAIMGNPKVKAHGKK VIN

Human  $\beta$  chain: AFSDGLAHLNLDNLKGTFA TLSELHCDKLHVDPENFRLLGN

Mouse  $\beta$  chain: AFNDGLKHLNLDNLKGTFA H LSELHCDKLHVDPENFRLLGN

Human  $\beta$  chain: VLVCVLAHHFGKEFTP PVQAAYQKVVAGVANALAHKYH

Mouse  $\beta$  chain: MI VI VLGHHLGKEFTP CAQAAFQKVVAGVASALAHKYH

There are a total of **27** mismatches, or  $(147 - 27) / 147 = 81.7\%$  identical

# Beta globins: Cont.

Human  $\beta$  chain: MVH **L** TPEEK**SAVT**ALWGKVN**DEVGGEALGRLL**

Chicken  $\beta$  chain: MVH**WTA**EEK**QL** I T**GL**WGKVN**AECGAEALARLL**

Human  $\beta$  chain: **VVYPWTQRFFESFGDLSTPDAVMGNPKVKAHGKKVLG**

Chicken  $\beta$  chain: **IVYPWTQRFFASFGNLSPTA I LGNPMVRAHGKKVLT**

Human  $\beta$  chain: **AFSDGLAHL**LDNLK**GT**F**ATL**SELHCDKLHVDPENF**RL**L**GN**

Chicken  $\beta$  chain: **SFGDAVKN**LDNIK **NTFSQ**LSELHCDKLHVDPENF**RL**L**GD**

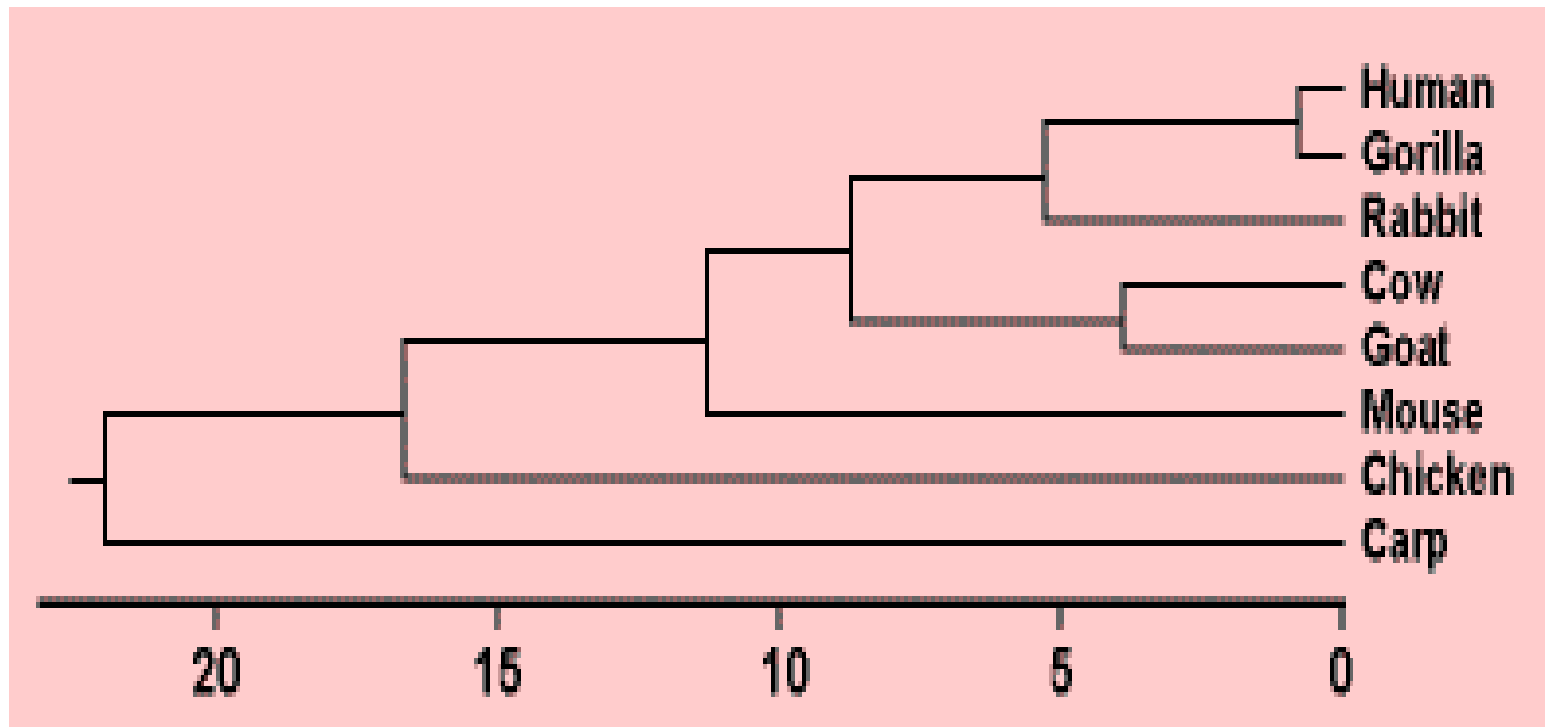
Human  $\beta$  chain: **VLVCVLAH**H**FGKE**FTP**PVQAAY** QK**VVAGVANALAHKYH**

Mouse  $\beta$  chain: **I L I I** VLA**AHF****SKD**FTP**ECQA****WQKLVRVVAHALARKYH**

-There are a total of **44** mismatches, or  $(147 - 44) / 147 = \mathbf{70.1\%}$  identical

- As expected, mouse  $\beta$  chain is '*closer*' to that of human than chicken's.

Molecular evolution can be visualized with phylogenetic tree.



Phylogenetic tree of Beta globin (Aligned using Clustal, PAM250)

---

# Mouse and Human overview

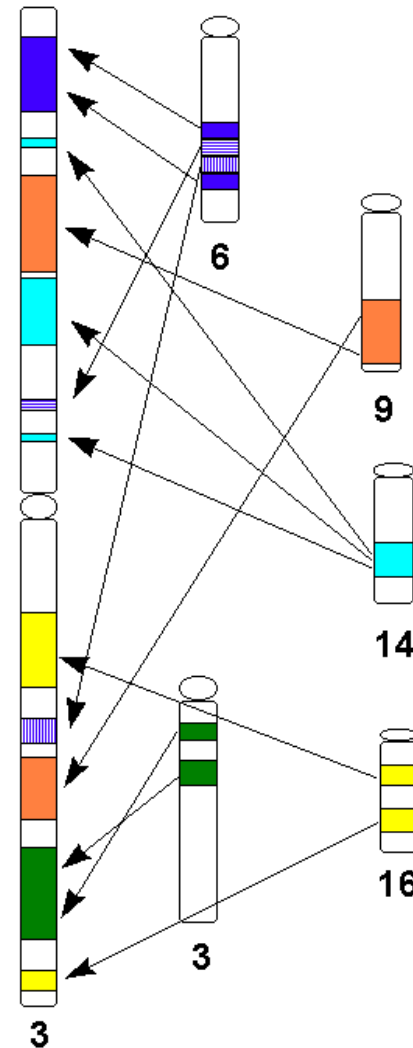
- Mouse has  $2.1 \times 10^9$  base pairs versus  $2.9 \times 10^9$  in human.
- About 95% of genetic material is shared.
- 99% of genes shared of about 30,000 total.
- The 300 genes that have no homologue in either species deal largely with immunity, detoxification, smell and sex\*

# Human and Mouse

Significant chromosomal rearranging occurred between the diverging point of humans and mice.

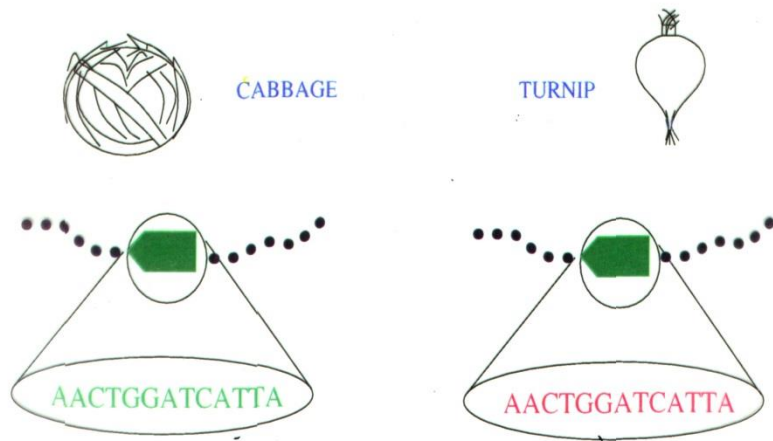
Here is a mapping of human chromosome 3.

It contains homologous sequences to at least 5 mouse chromosomes.



# Important discovery

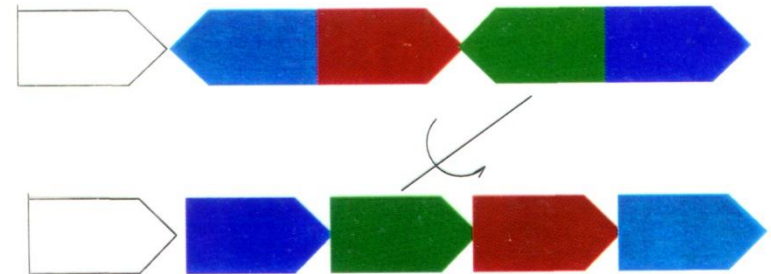
## GENE SEQUENCE COMPARISON



AACTGGATCATT  
AACTGGATCATT

Comparing gene sequences yields  
no evolutionary information

## GENE ORDER COMPARISON

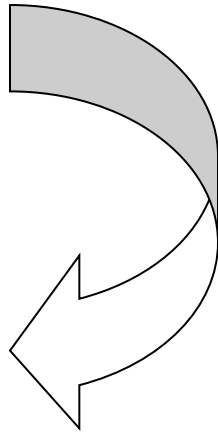


Evolution is manifested as the  
divergence in Gene Order

# DNA Reversal



Break  
and  
Invert



5' ATG CCTGTA CTA 3'

3' TAC GGACATGAT 5'



5' ATG TACAGG CTA 3'

3' TAC ATGTCCGAT 5'

---

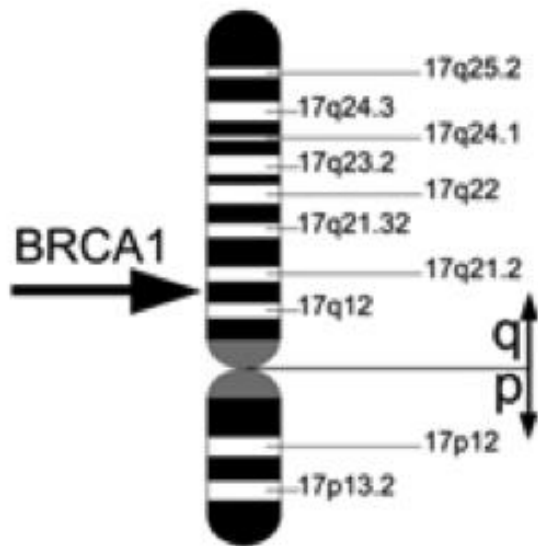
# Waardenburg's syndrome

- Genetic disorder
- Characterized by loss of hearing and pigmentary dysphasia
- Found on human chromosome 2

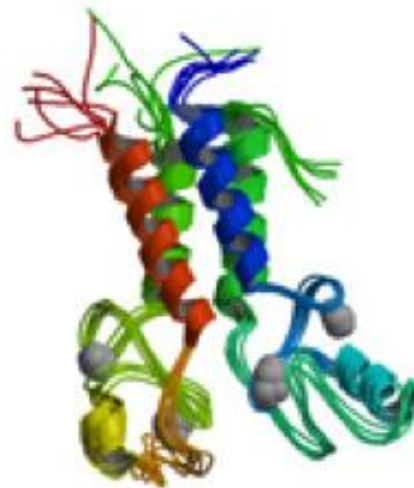




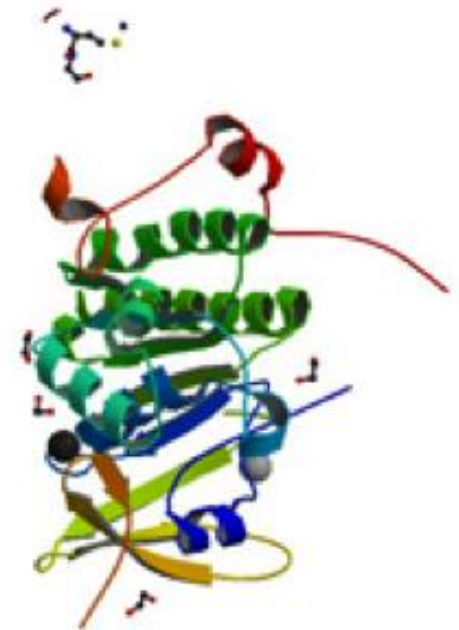
## Chromosome 17



## BRCA1



## BRCA2



---

# It is Sequenced, What's Next?

- **Tracing Phylogeny**
    - Finding family relationships between species by tracking similarities between species.
  - **Gene Annotation (cooperative genomics)**
    - Comparison of similar species.
  - **Proteomics**
    - From DNA sequence to a folded protein.
  - **Determining Regulatory Networks**
    - How the genes react to certain stimuli?
    - How cancer progresses?
-

---

# Sequence Driven Problems in Bioinformatics

- Genomics
    - Fragment assembly of the DNA sequence.
      - Not possible to read entire sequence.
      - Cut up into small fragments using restriction enzymes.
      - Then need to do fragment assembly. Overlapping similarities to matching fragments.
      - N-P complete problem.
    - Finding Genes
      - Identify open reading frames
        - Exons are spliced out.
        - Junk in between genes
-

---

# *More Computing problems...*

- Proteomics
    - Protein Folding
      - 1D Sequence → 3D Structure
      - What drives this process?
    - Identification of functional domains in protein's sequence
      - Determining functional pieces in proteins.
  - *Most problems need mathematical solutions & For Large data they need software*
-

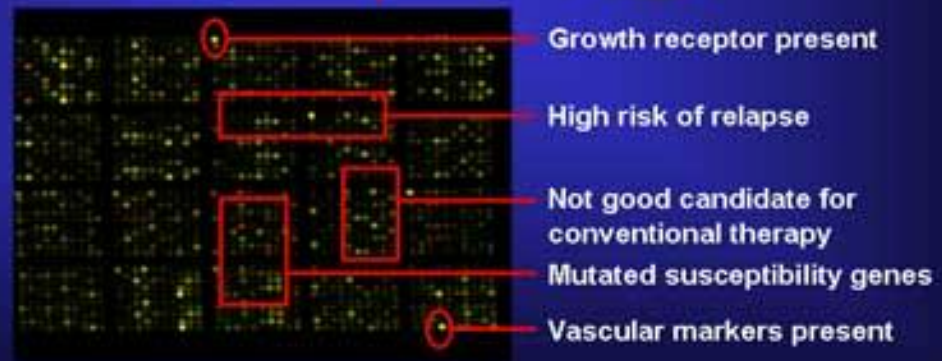
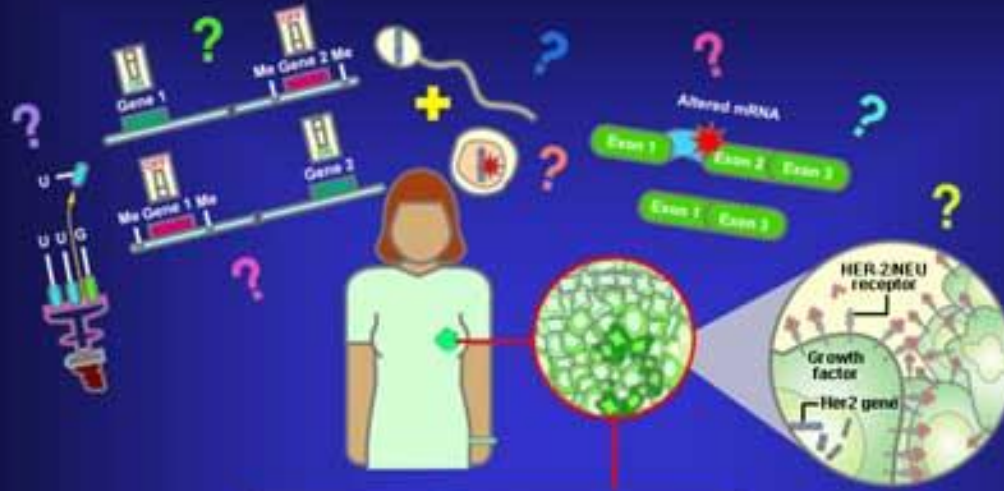
---

# Biological Databases

- Online databases to archive, search, and run programs on – massive amount of data
  - **NCBI GeneBank** <http://ncbi.nih.gov>  
Huge collection of databases, the most prominent being the nucleotide sequence database
  - **Protein Data Bank** <http://www.pdb.org>  
Database of protein 3D-structures
  - **SWISSPROT** <http://www.expasy.org/sprot/>
    - Database of annotated protein sequences
  - **PROSITE** <http://kr.expasy.org/prosite>  
Database of protein active site motifs
-

# How gene expressions translate to Prognosis or Diagnosis: A Cancer Question

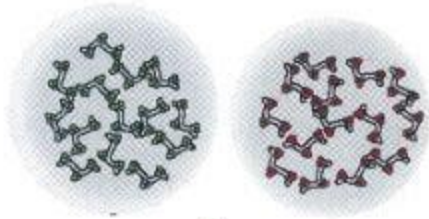
## A Daunting Challenge



Adapted by Joanna Kelly, © 2004

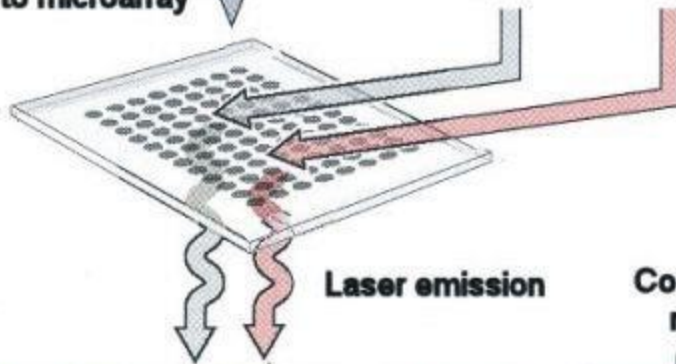
Make cDNA reverse transcript  
Label cDNAs with fluorescent dyes

Control      Experimental



Hybridization  
to microarray

Laser excitation  
at dye-specific Hz



Red = "up-regulation"

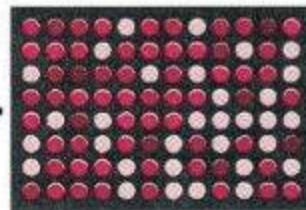
Green = "down-regulation"

Black = constitutive  
expression

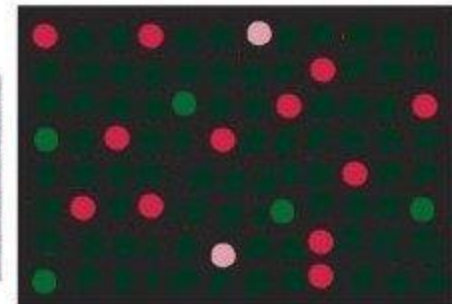
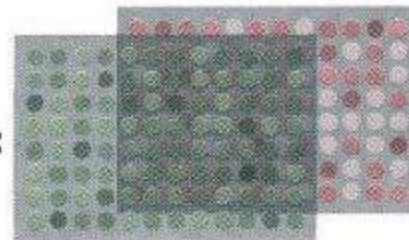
Computer calculates  
ratio of intensity



+



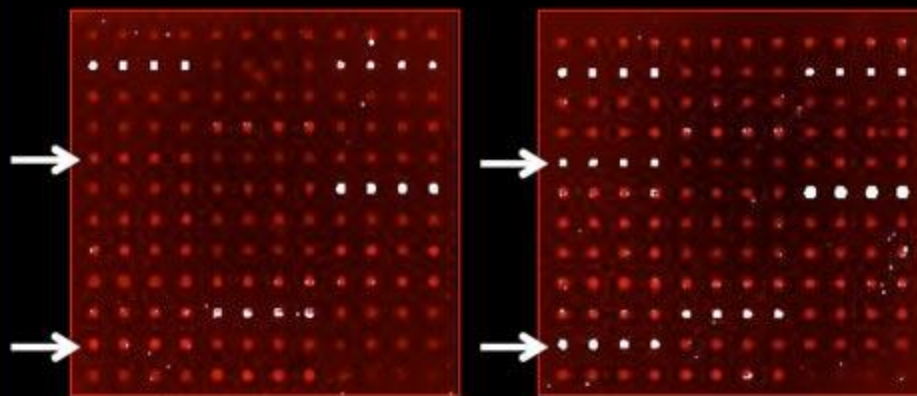
=



# Patient Serum Samples

Normal

Prostate  
Cancer

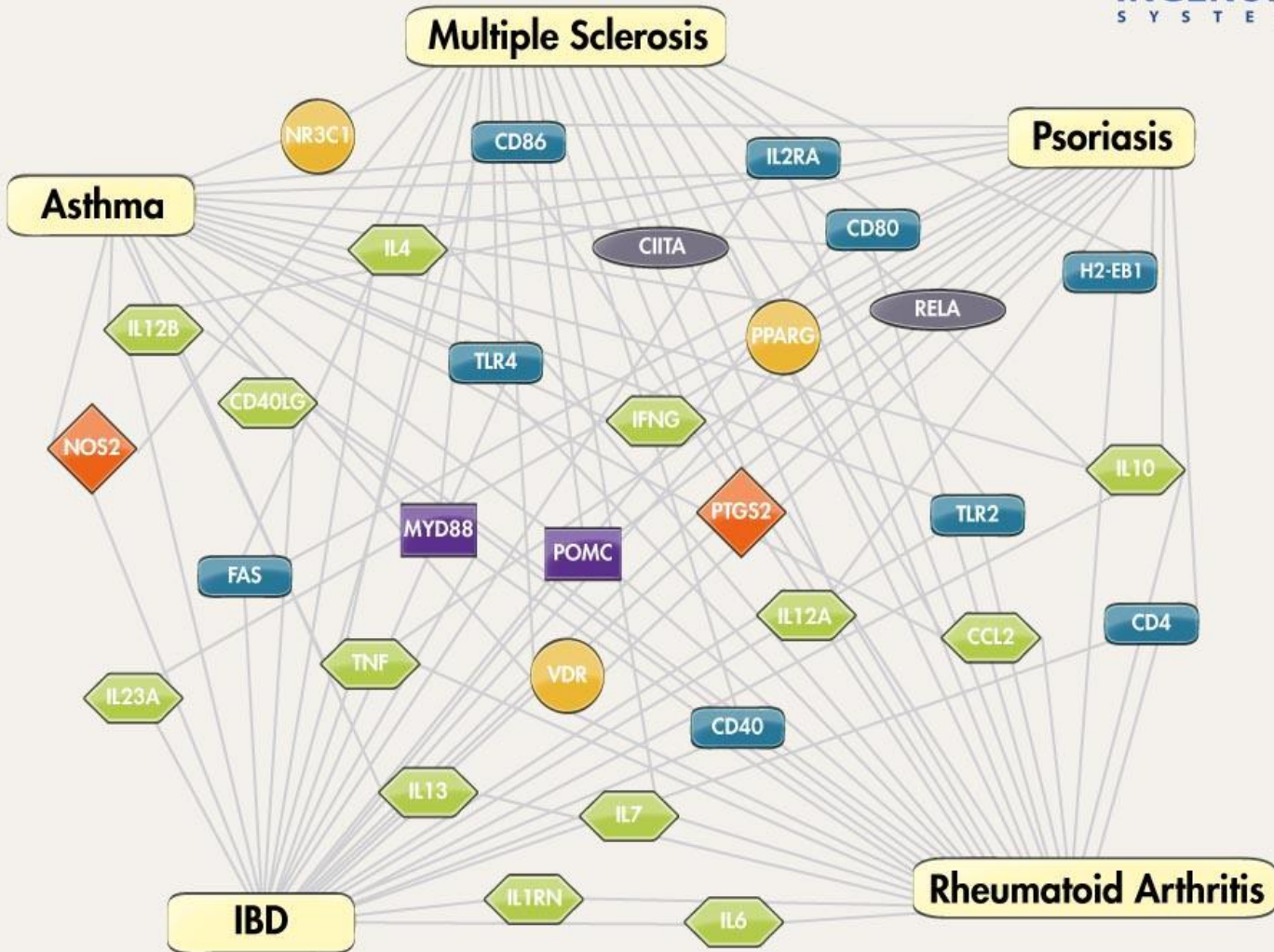


GAT™ polyclonal antibodies spotted in quadruplicate



# Microarray Experiments lead to understanding Gene Regulatory Networks

Powered By  
**INGENUITY**  
SYSTEMS



# Big Data?



"Your recent Amazon purchases, Tweet score and location history makes you 23.5% welcome here."

---

# Current Thoughts on Big-data & Biology – a Conference

<http://www.triconference.com/Bioinformatics-Genome/>

***“Bioinformatics for Big Data:  
How Applications of Big Data will Drive Research  
Forward”***

**February 10-12, 2014 | Moscone North Convention Center | San Francisco, CA**

---

---

# Session: TECHNOLOGIES GENERATING BIG BIOMEDICAL DATA

## **11:00 Cancer Genomics**

*David Haussler, Ph.D., Distinguished Professor and Director,  
Center for Biomolecular Science & Engineering, Univ. of California Santa Cruz*

UCSC has built the Cancer Genomics Hub (CGHub) for the US National Cancer Institute, designed to hold up to 5 petabytes of research genomics data (up to 50,000 whole genomes), including data for all major NCI projects. To date it has served more than 8.3 petabytes of data to more than 300 research labs. Cancer is exceedingly complex, with thousands of subtypes involving an immense number of different combinations of mutations. The only way we will understand it is to gather together DNA data from many thousands of cancer genomes so that we have the statistical power to distinguish between recurring combinations of mutations that drive cancer progression and "passenger" mutations that occur by random chance. Currently, with the exception of a few projects such as ICGC and TCGA, most cancer genomics research is taking place in research silos, with little opportunity for data sharing. If this trend continues, we lose an incredible opportunity. Soon cancer genome sequencing will be widespread in clinical practice, making it possible in principle to study as many as a million cancer genomes. For these data to also have impact on understanding cancer, we must begin soon to move data into a global cloud storage and computing system, and design mechanisms that allow clinical data to be used in research with appropriate patient consent. A global alliance for sharing genomic and clinical data is emerging to address this problem. This is an opportunity we cannot turn away from, but involves both social and technical challenges.

***Translation:  
Huge database of mutated gene's***

---

---

# Session: DATA STORAGE AND MAINTENANCE

## **2:20 Implementing Big Data Analysis and Archival Solutions for NGS Data**

*Zhiyan Fu, Ph.D., Chief Scientific Computing Officer, Genome Institute of Singapore (A\*STAR)*

This presentation shows the latest development in big data analysis, compression and storage management. It provides a practical case to implement the big data technologies to a mid-size genome center. Attendees will understand the challenges of big data life-cycle management in a genome center and see how the latest big data technologies are implemented, and the pros and cons of some of the techniques, including Hadoop, HDF5, and different NGS compression algorithms evaluated by GIS.

***Translation:***

***Experience in managing Next Gen. Sequencing Data storage***

---

---

# Session: HOW BIG DATA WILL DRIVE RESEARCH FORWARD

## **1:45 Harnessing Big Data to Accelerate Drug Development**

*Vinod Kumar, Ph.D., Senior Investigator, Computational Biology, GlaxoSmithKline Pharmaceuticals*

With the rapid development of high-throughput technologies and ever-increasing accumulation of whole genome-level datasets, an increasing number of diseases and drugs can be comprehensively characterized by the changes they induce in gene expression, protein, metabolites and phenotypes. Integrating and querying such large volumes of data, often spanning domains and residing in diverse sources, constitutes a significant obstacle. This talk presents two distinct approaches that utilize these data types to systematically evaluate and suggest new disease indications for new and existing drugs.

***Translation:***

***Predict new disease indications from Data for Drug Modeling, first by software***

---

---

# Session: BIG DATA DRIVING PERSONALIZED MEDICINE

**5:05 It's Not Just About Big Data...**

**Big Analytics for Identifying What Works and for Whom in Healthcare**

*Iya Khalil, Ph.D., Executive Vice President and Co-Founder, GNS Healthcare*

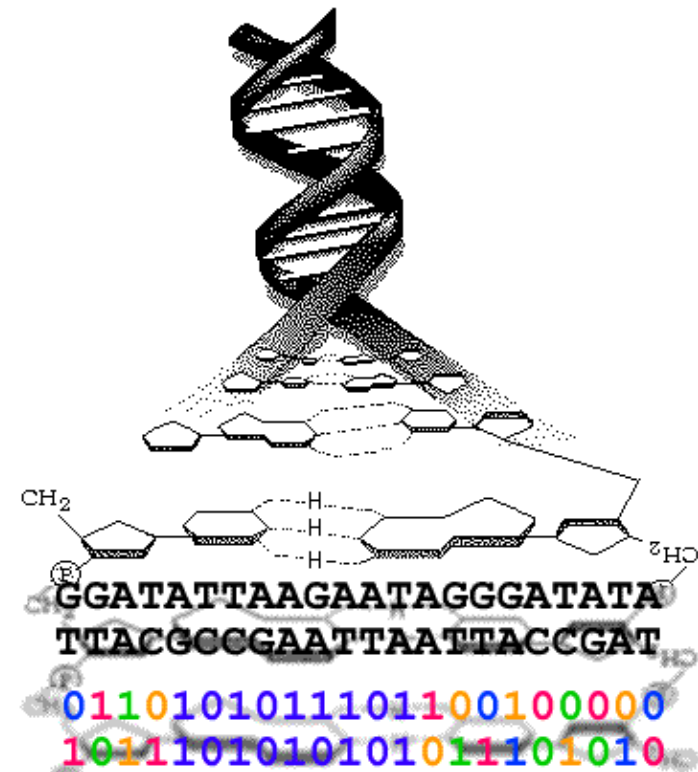
We are living in the era of big data in healthcare, with unprecedented ability to collect data at multiple levels (genomic/'omic', phenotypic, health records, mobile health, etc.) and at scale. The key will be leveraging advanced analytics and appropriate feedback loops to identify what works on an individual patient level.

***Translation:***

***Smart Algorithms needed for connecting multi-omics Big  
Data***

---

*Thank you!*



Debasis Mitra

Professor, Computer Science, FIT

[dmitra@cs.fit.edu](mailto:dmitra@cs.fit.edu)

