



# CSCI 1810

## Computational Molecular Biology

Of Sea Urchins, Birds and Humans



# Overview

---

1. The Human Genome
2. The Molecular Biology Dogma: DNA, RNA and Protein
3. **Beautiful Algorithms: Rigorous, Practical, Elegant code**
4. Chapter 1: Sequence Alignment Algorithms
5. Chapter 2: Combinatorial Pattern Matching Algorithms
6. Chapter 3: Phylogenetic Trees Algorithms
7. Chapter 4: Machine Learning Methods: Hidden Markov Models Algorithms
8. Chapter 5: Genome Assembly Algorithms (Introduction)
9. Chapter 6: Genomic Privacy (Introduction)
10. The Bioinformatician as a Detective – two puzzles:
  - The Adventures of the Dancing Men code, by Sherlock Holmes/Arthur Conan Doyle
  - The Prison code, a code used in a prison in California

# Beautiful Algorithms

---

- **Rigorous:** state-of-the-art, mathematical analysis of their accuracy
- **Practical:** very efficient, work on large data sets
- **Elegant code:** “simplicity is the ultimate sophistication”
- **von Neumann’s “esthetic criteria”**  
many applications to different areas



- 
- John von Neumann’s **“Beautiful” criteria**
  - *“By a model is meant a mathematical construct which, with the addition of certain verbal interpretations, describes the observed phenomena... Furthermore, it must satisfy certain esthetic criteria – that is, in relation to how much it describes, it must be rather simple ...One cannot tell exactly how “simple” simple is. ...Simplicity is largely a matter of historical background, of previous conditioning, of antecedents, of customary procedures, and it is very much a function of what is explained by it.”*
  - – John von Neumann



---

# Evolution

# Evolution

---



Theodosius Dobzhansky (1900-1975)

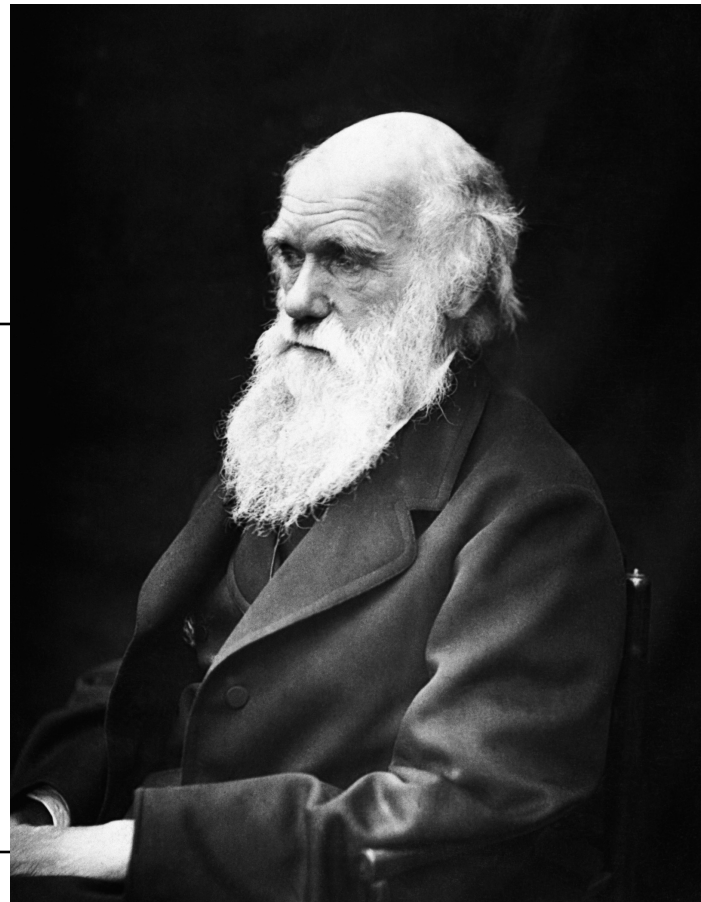
**Nothing in Biology Makes Sense  
Except in the Light of Evolution**



# Darwin's Finches



and Coco

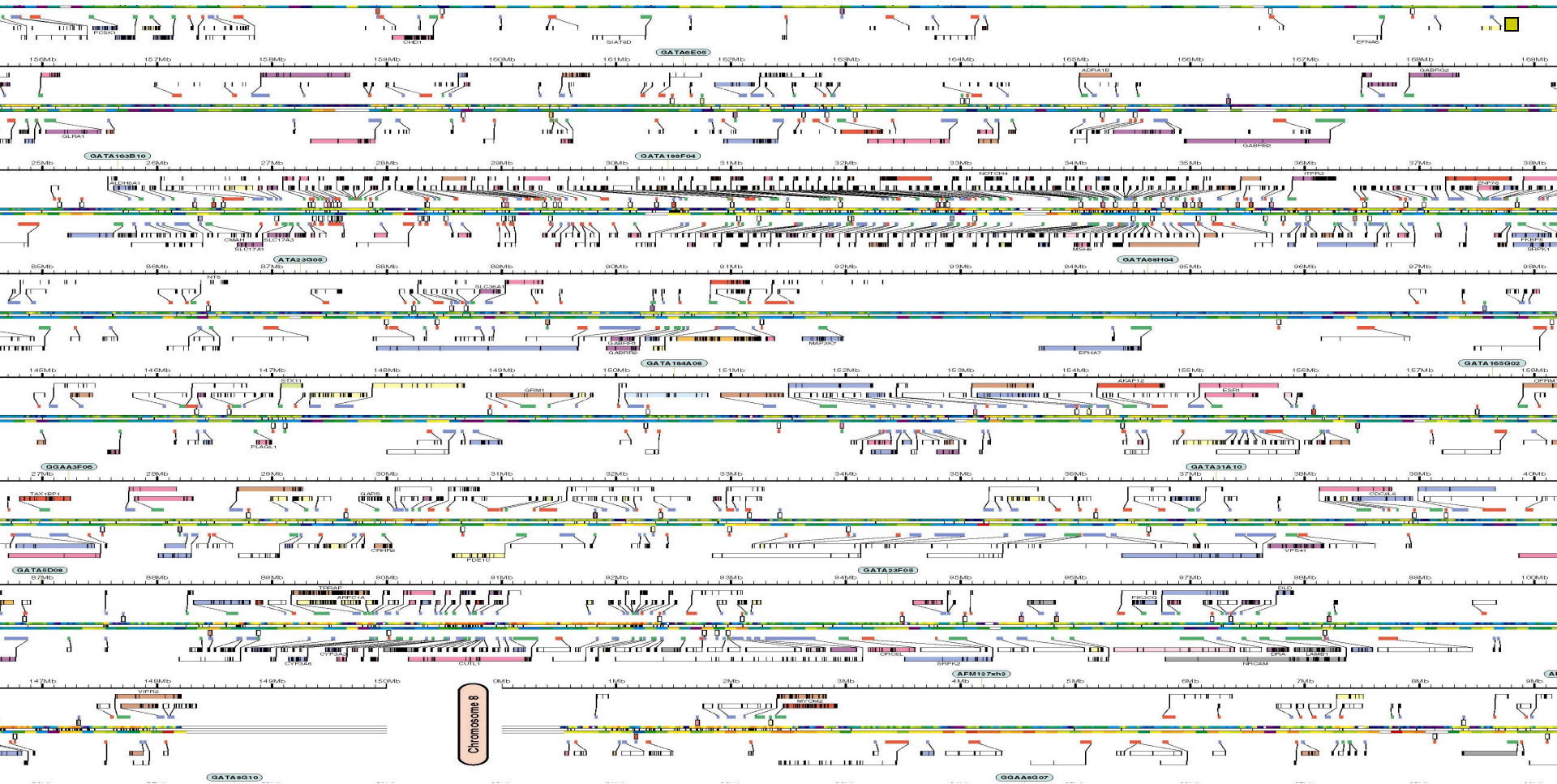




---

# The Genome

# The Sequence of the Human Genome





# Science

16 February 2001

Vol. 291 No. 5507  
Pages 1145-1434 \$9

## THE HUMAN GENOME



AMERICAN ASSOCIATION FOR THE ADVANCEMENT OF SCIENCE



# Biomolecular Data

```
A G C T T A C T A A T C C G G G C C G A A T T A G G T C
A G T T T A T T A A T T C G A G C T G A A C T A G G T C
A G T C T A T T A A T T C G A G C A G A A C T T G G T C
A G T T T A T T A A T T C G A G C T G A A C T T G G C C
A G T C T A C T A A T T C G A G C T G A A T T A G G T C
A G A T T A T T A A T T C G A G C T G A A C T T G G T C
A G A T T G C T A A T T C G A G C C G A A T T A G G T C
A G A T T A T T A A T C C G G G C T G A A T T A G G T C
A G T C T A T T A A T T C G A G C T G A A T T A G G A C
A G C T T A T T A A T T C G T G C T G A A C T C G G A C
A G C T T A T T A A T T C G A G C T G A A C T C G G A C
A G C T T A T T A A T T C G A G C C G A A C T C G G G C
A G T C T T T T A A T T C G A G C T G A A T T A G G A C
```

# Biomolecular Data

---

“Data! Data! Data!” he cried impatiently. “I can’t make bricks without clay.”

Sherlock Holmes



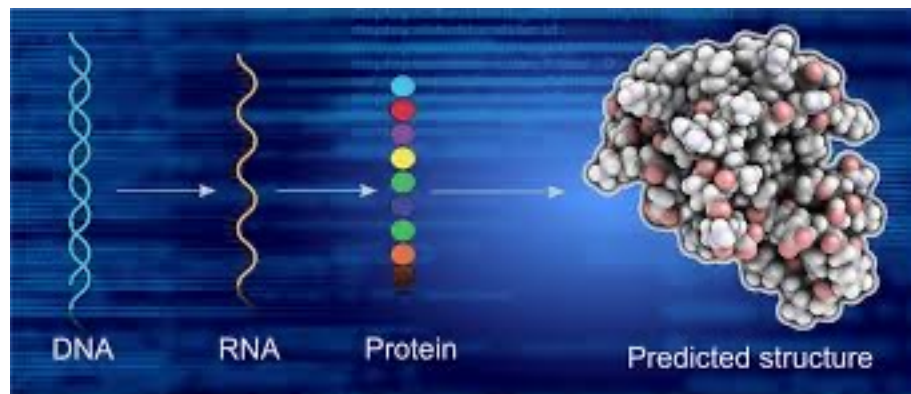
“The more I see the less I know for sure.”

John Lennon



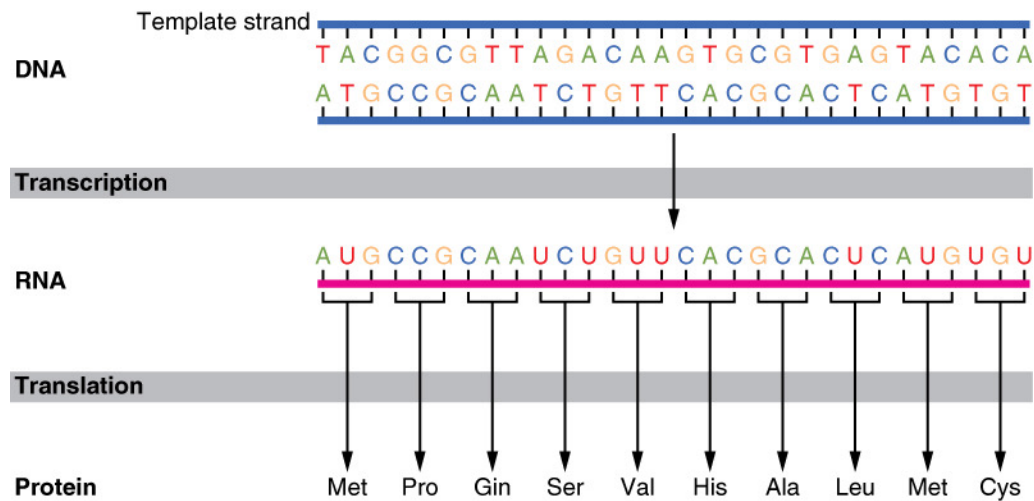


# The Dogma of Molecular Biology



# The Molecular Biology Dogma

---



# Transcription

gene

Exon 1 Intron 1 Exon 2 Intron 2 Exon 3

**DNA**

ACGTCT	GTACTGCATT	AGCGATG	CATACG	ATGCATGCAA	GGCATA
TGCAGAT	CATGACGTA	TCGCTAC	GTATGC	TACGTACGTTT	CCGTATG



RNA polymerase

nuclear factors

RNA

ACGTCT	GTACTGCATT	AGCGATG	CATACG	ATGCATGCAA	GGCATA
TGCAGAT	CATGACGTA	TCGCTAC	GTATGC	TACGTACGTTT	CCGTATG

GUAC



**hRNA**  
(heteronuclear)

splicing

**mRNA**  
(messenger)

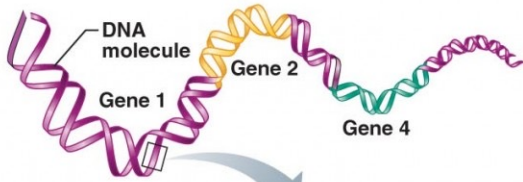
polyA tail

GUACUGCAUU	AGCGAUG	CAUACG	AUGCAUGCAA	GGCAUAC
------------	---------	--------	------------	---------

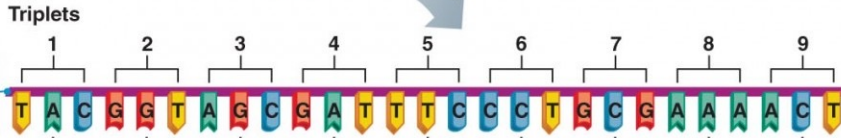
  

GUACUGCAUU	CAUACG	GGCAUAC	AAAAAAAAAAAA
------------	--------	---------	--------------

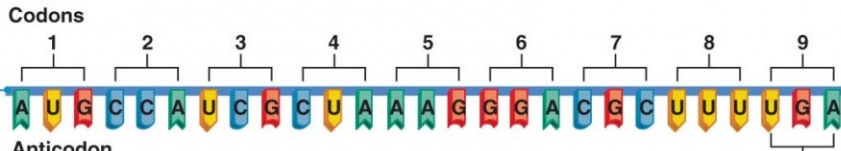




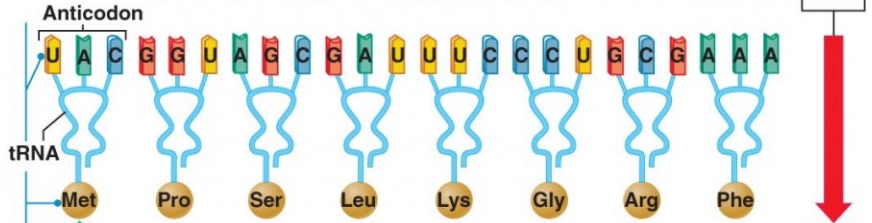
**DNA:** DNA base sequence (triplets) of the gene codes for synthesis of a particular polypeptide chain



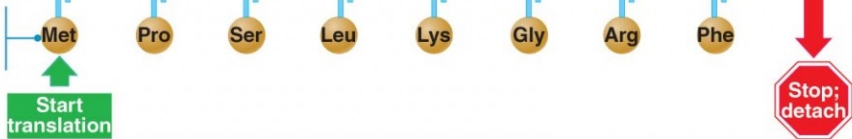
**mRNA:** Base sequence (codons) of the transcribed mRNA



**tRNA:** Consecutive base sequences of tRNA anticodons recognize the mRNA codons calling for the amino acids they transport

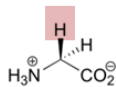


**Polypeptide:** Amino acid sequence of the polypeptide chain

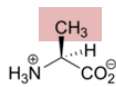


# The set of 20 amino acids

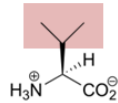
## Nonpolar, aliphatic side groups



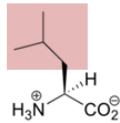
Glycine  
Gly, G



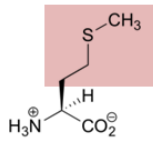
Alanine  
Ala, A



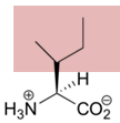
Valine  
Val, V



Leucine  
Leu, L

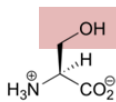


Methionine  
Met, M

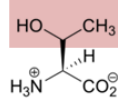


Isoleucine  
Ile, I

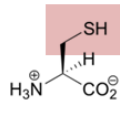
## Polar, uncharged side groups



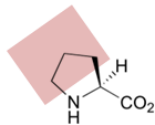
Serine  
Ser, S



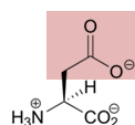
Threonine  
Thr, T



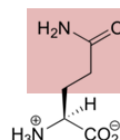
Cysteine  
Cys, C



Proline  
Pro, P

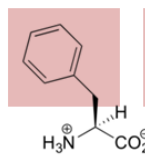


Aspartate  
Asp, D

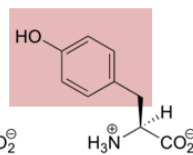


Glutamine  
Gln, Q

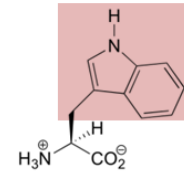
## Aromatic side groups



Phenylalanine  
Phe, F

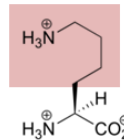


Tyrosine  
Tyr, Y

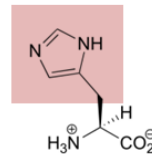


Tryptophan  
Trp, W

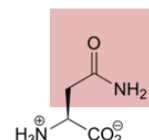
## Positively charged side groups



Lysine  
Lys, K

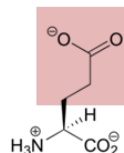


Histidine  
His, H

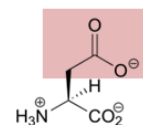


Asparagine  
Asn, N

## Negatively charged side groups



Glutamate  
Glu, E

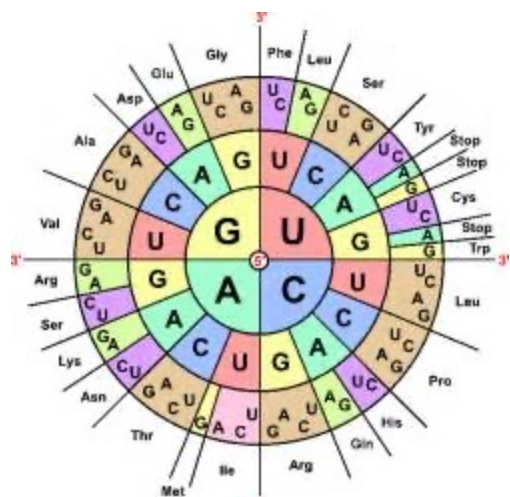


Aspartate  
Asp, D

# The Genetic Code

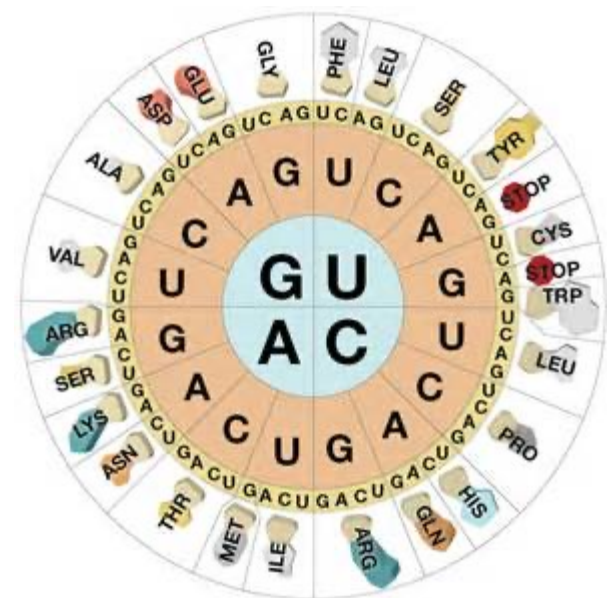
		Second base				Third base
		U	C	A	G	
First base	U	UUU - Phenylalanine UUC - Phenylalanine UUA - Leucine UUG - Leucine	UCU - Serine UCC - Serine UCA - Serine UCG - Serine	UAU - Tyrosine UAC - Tyrosine UAA - Stop codon UAG - Stop codon	UGU - Cysteine UGC - Cysteine UGA - Stop codon UGG - Tryptophan	U C A G
	C	CUU - Leucine CUC - Leucine CUA - Leucine CUG - Leucine	CCU - Proline CCC - Proline CCA - Proline CCG - Proline	CAU - Histidine CAC - Histidine CAA - Glutamine CAG - Glutamine	CGU - Arginine CGC - Arginine CGA - Arginine CGG - Arginine	U C A G
	A	AUU - Isoleucine AUC - Isoleucine AUA - Methionine start codon AUG - Methionine start codon	ACU - Threonine ACC - Threonine ACA - Threonine ACG - Threonine	AAU - Asparagine AAC - Asparagine AAA - Lysine AAG - Lysine	AGU - Serine AGC - Serine AGA - Arginine AGG - Arginine	U C A G
G	GUU - Valine GUC - Valine GUA - Valine GUG - Valine	GCU - Alanine GCC - Alanine GCA - Alanine GCG - Alanine	GAU - Aspartic acid GAC - Aspartic acid GAA - Glutamic acid GAG - Glutamic acid	GGU - Glycine GGC - Glycine GGA - Glycine GGG - Glycine	U C A G	

		Second Letter				3rd letter
		U	C	A	G	
1st letter	U	UUU - Phe UUC - Phe UUA - Leu UUG - Leu	UCU - Ser UCC - Ser UCA - Ser UCG - Ser	UAU - Tyr UAC - Tyr UAA - Stop UAG - Stop	UGU - Cys UGC - Cys UGA - Stop UGG - Trp	U C A G
	C	CUU - Leu CUC - Leu CUA - Leu CUG - Leu	CCU - Pro CCC - Pro CCA - Pro CCG - Pro	CAU - His CAC - His CAA - Gln CAG - Gln	CGU - Arg CGC - Arg CGA - Arg CGG - Arg	U C A G
	A	AUU - Ile AUC - Ile AUA - Met AUG - Met	ACU - Thr ACC - Thr ACA - Thr ACG - Thr	AAU - Asn AAC - Asn AAA - Lys AAG - Lys	AGU - Ser AGC - Ser AGA - Arg AGG - Arg	U C A G
	G	GUU - Val GUC - Val GUA - Val GUG - Val	GCU - Ala GCC - Ala GCA - Ala GCG - Ala	GAU - Asp GAC - Asp GAA - Glu GAG - Glu	GGU - Gly GGC - Gly GGA - Gly GGG - Gly	U C A G



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

  Hydrophobic - Imino   
   Hydrophobic - Aliphatic   
   Polar - Neutral  
  Hydrophobic - Aromatic   
   Polar - Acid   
   Polar - Basic





---

# Genetic Variation

**SNPs & HAPLOTYPES**

# Single Nucleotide Polymorphism (SNP)

---

GATTAGATC**G**CGATAGAG  
GATTAGATC**T**CGATAGAG

A SNP is a position in a genome at which two or more different bases occur in the population, each with a frequency  $>1\%$ .

The two alleles at the site are **G** and **T**

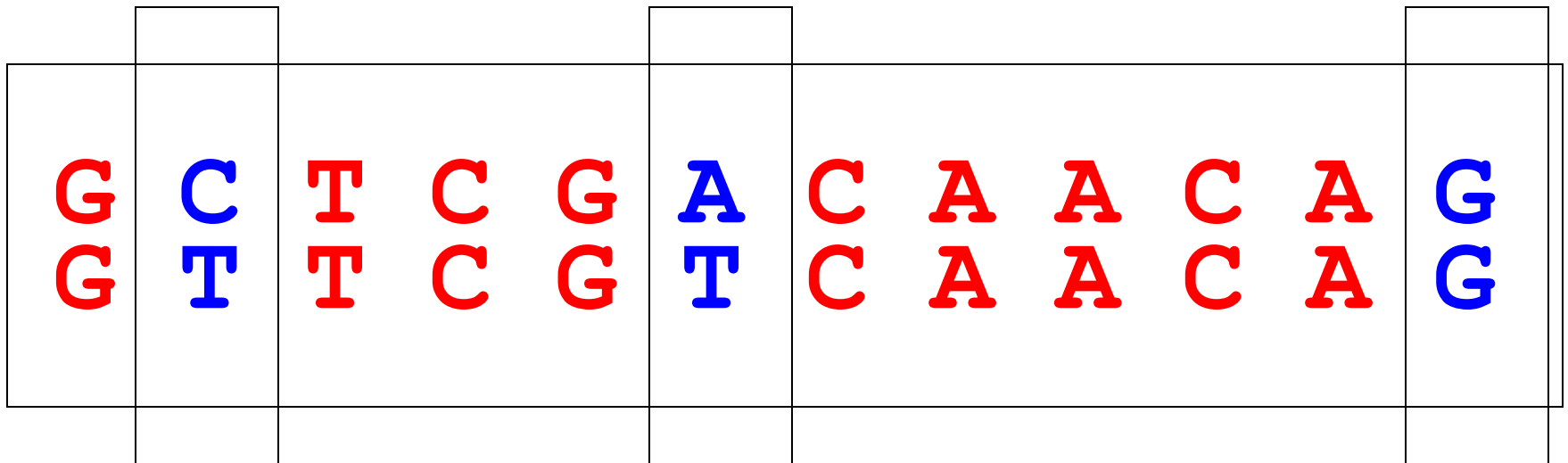
- The most abundant type of polymorphism



# Haplotype

C A G  
T T G

Haplotypes



SNP

SNP

SNP

Two individuals

# Mutations

```
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
ATCTATATGCGTACGTACGTACGTAC
```

Infinite Sites Assumption:  
Each site mutates at most once



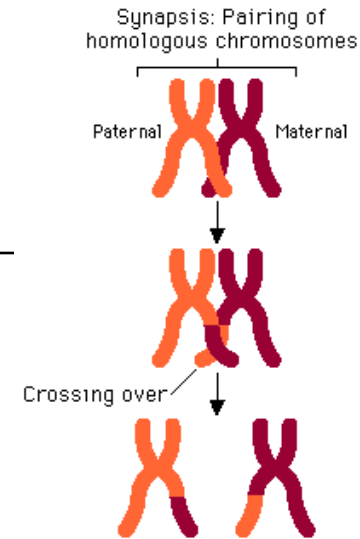
# Haplotype Pattern

---

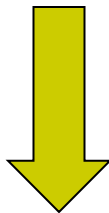
C	A	G	T	0	0	0	0
T	T	G	A	1	1	0	1
C	A	T	G	0	0	1	0
C	T	G	T	0	1	0	1

At each SNP site label the two alleles as 0 and 1.  
The choice which allele is 0 and which one is 1  
is arbitrary.

# Recombination



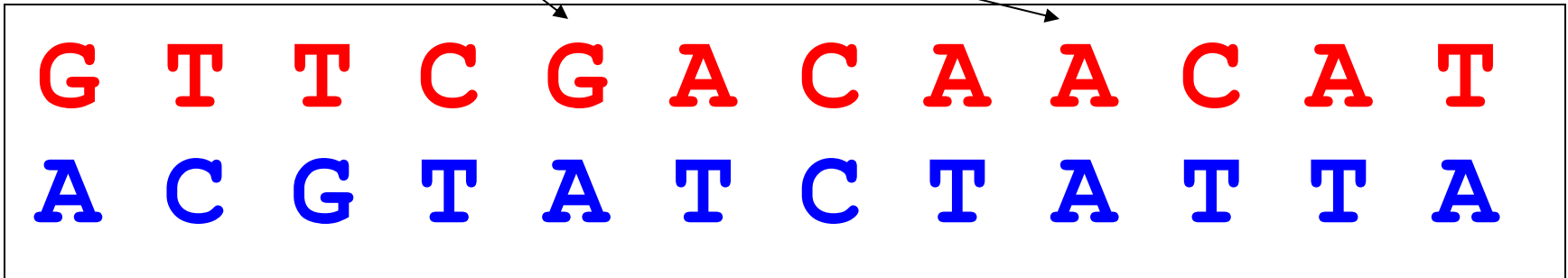
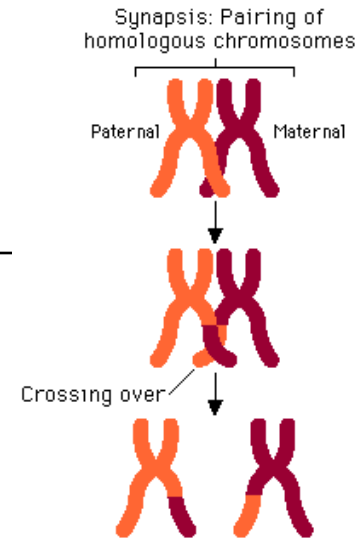
G	T	T	C	G	A	C	A	A	C	A	T
A	C	G	T	A	T	C	T	A	T	T	A



G T T C G A C T A T T A

# Recombination

The two alleles are linked, I.e., they are “**traveling together**”



Recombination  
disrupts the linkage

?

G T T C G A C T A T T A

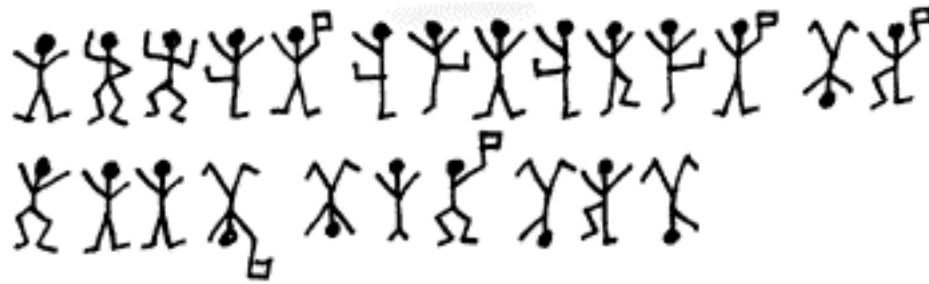
CACAGCCTGGATAACAGGAGGACCTTGATGCTCCTGGCACAAATGAGCAGAATCTCT  
CCTTCCTCCTGTCTGATGGACAGACATGACTTTGGATTCCCCAGGAGGAGTTTGAT  
GGCAACCAGTTCCAGAAGGCTCCAGCCATCTCTGTCCTCCATGAGCTGATCCAGCAG  
ATCTTCAACCTCTTTACCACAAAAGATTCATCTGCTGCTTGGGATGAGGACCTCCTA  
GACAAATTCTGCACCGAACTCTACCAGCAGCTGAATGACTTGGGAAGCCTGTGTGATG  
CAGGAGGAGAGGGTGGGAGAACTCCCCTGATGAATGCGGACTCCATCTTGGCTGTG  
AAGAAATACTTCGAAGAATCACTCTCTATCTGACAGAGAAGAAATACAGCCCTTGT  
GCCTGGGAGGTTGTCAGAGCAGAAATCATGAGATCCTCTCTTTATCAACAAACTTGC  
AAGAAAGATTAAGGAGGAAGGAATAA, TGTGATCTCCCTGAGACCCACAGCCTGGA  
TAACAGGAGGACCTTGATGCTCCTGGCACAAATGAGCAGAATCTCTCCTTCCTCCTG  
TCTGATGGACAGACATGACTTTGGATTCCCCAGGAGGAGTTTGATGGCAACCAGTT  
CCAGAAGGCTCCAGCCATCTCTGTCCTCCATGAGCTGATCCAGCAGATCTTCAACCT

What is the meaning of this DNA sequence?

A code to break!

# Can you break this code?

---



# Chapter 1: Sequence Alignment

---





Avrilla Xinyue Qian

# Chapter 1: Sequence Alignment Algorithms

## Local Alignment

### Pairwise Sequence Alignment

Target Sequence

5' ACTACTAGATTACTTACGGATCAGGTACTTTAGAGGCTTGCAACCA 3'

||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||

Query Sequence 5' TACTCACGGATGAGGTACTTTAGAGGC 3'

## Global Alignment

Target Sequence

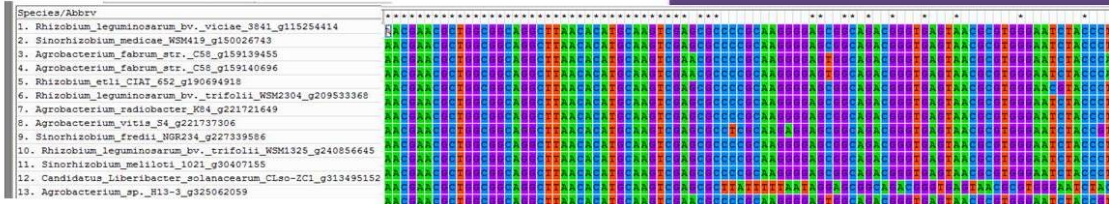
5' ACTACTAGATTACTTACGGATCAGGTACTTTAGAGGCTTGCAACCA 3'

||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||

5' ACTACTAGATT----ACGGATC--GTACTTTAGAGGCTAGCAACCA 3'

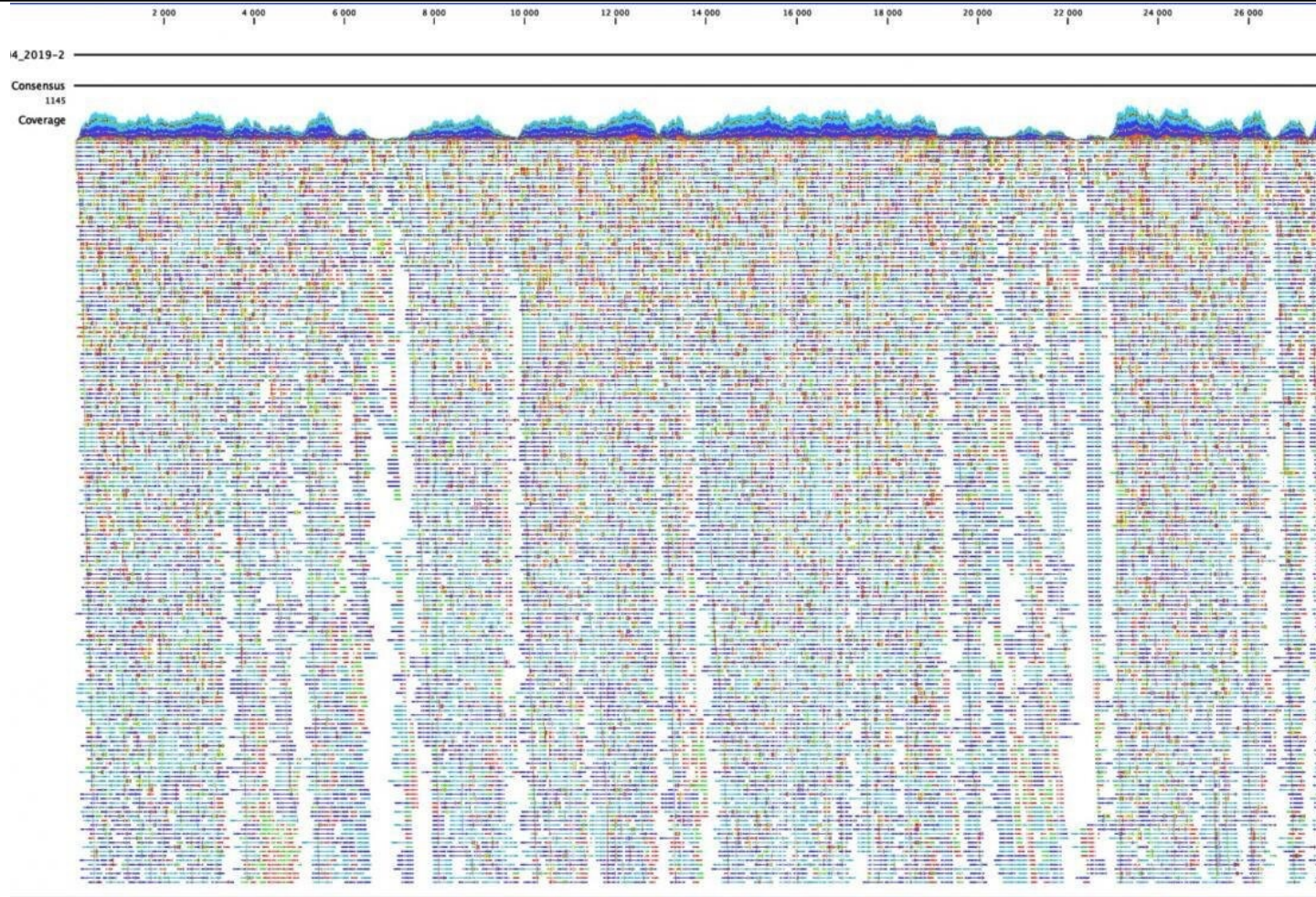
Query Sequence

### Multiple Sequence Alignment (MSA)









Whole genome sequence of the 2019-nCoV **coronavirus**, in one of the first French cases, made at the Institut Pasteur (Paris), using a unique Platform (P2M), open to all French National Reference Centers. Credit: Institut Pasteur/CNR of respiratory infection viruses.



# Margaret Dayhoff & PAM Similarity Matrices

---





# **Dr. Margaret Oakley Dayhoff**

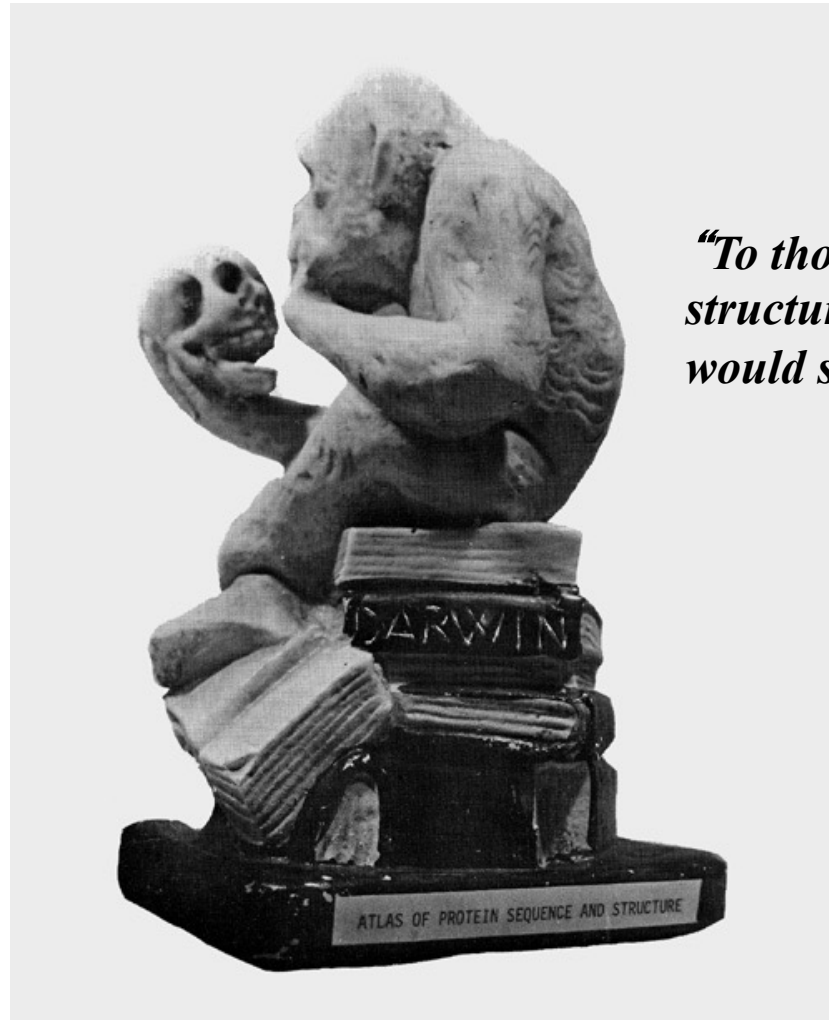
## **The Mother & Father of Bioinformatics**

---





## The Atlas of Protein Sequence and Structure 1972



*“To those who would know the biochemical structure, function and origin of man and would strive to improve his lot.”*

Group & Subgroup Names	Amino Acid Residue	Group Properties
<b>Hydrophilic</b> -Small Aliphatic  -Acid amide  -Acid  -Hydroxyl	Alanine Proline Glycine  Glutamine Asparagine  Glutamic Acid Aspartic Acid  Serine Threonine	<b>Small, Simple, Hydrophilic</b> Not hydrophobic, smallest  Slightly basic, amide, carbonyl  Acid, carbonyl  Hydroxyl, small
<b>Sulfhydryl</b>	Cysteine	<b>Uniquely Reactive, Small</b>
<b>Aliphatic</b>	Valine Isoleucine Methionine Leucine	<b>Hydrophobic</b> Similarly branched
<b>Basic</b>	Lysine Arginine Histidine	<b>Basic, Nitrogen, Large</b>
<b>Aromatic</b>	Phenylalanine Tyrosine Tryptophan	<b>Aromatic Rings, Hydrophobic, Large</b>
<b>Special</b>	Histidine Tryptophan  Cysteine Serine  Phenylalanine Leucine Isoleucine Methionine	Heterocyclic rings  Close similarity in shape  Hydrophobic; similar size





# The Smith Waterman Algorithm

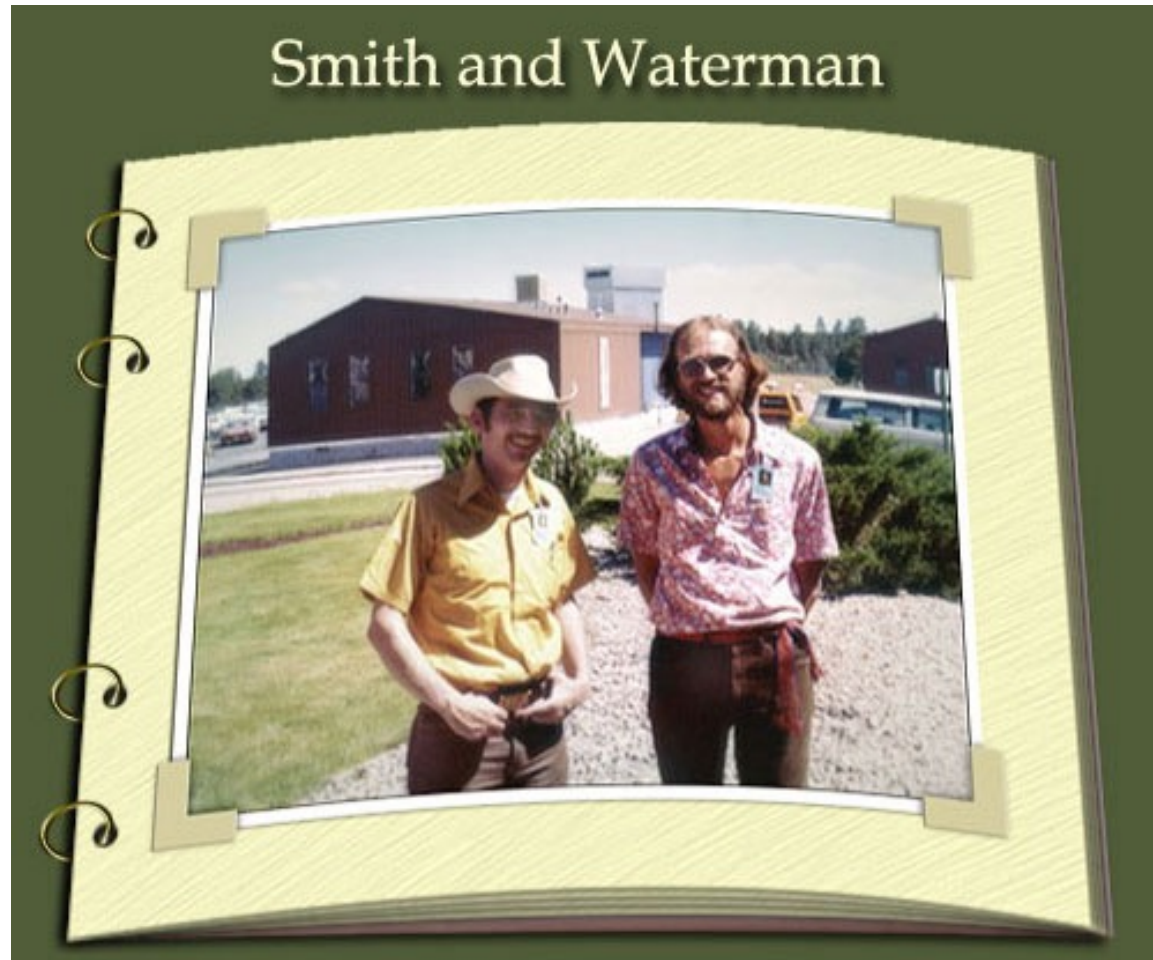
---



# Smith and Waterman at Los Alamos, New Mexico

Photo by David Lipman, Taken Summer of 1980

---





# Viral src gene products are related to the catalytic chain of mammalian cAMP-dependent protein kinase.

---

## **AUTHORS**

W. C. Barker and M. O. Dayhoff

## **ABSTRACT**

The transforming protein sequences translated from the Rous avian and Moloney murine sarcoma virus src genes are shown to be related to the catalytic chain of bovine cAMP-dependent protein kinase (ATP:protein phosphotransferase, EC 2.7.1.37). The avian transforming protein, also a protein kinase, shows greatest homology with the bovine protein kinase in the carboxyl-terminal half, where the protein kinase activity is localized. Moreover, lysine occurs in the inferred transforming protein sequences at the position homologous with the proposed ATP-binding lysine of the bovine protein kinase. This relationship is consistent with the hypothesis that the src genes originated in the host genomes, in which they are members of a superfamily of distantly related protein kinases that are normal constituents of mammalian cells. In the host, these sequences are much more highly conserved than in the viruses.



# Viral *src* gene products are related to the catalytic chain of mammalian cAMP-dependent protein kinase

---

1 BOV-PK QIEHTL NEKRI - -LQAV NFPF LVKLEFSF KDNSNLYMVM EYV PGGE MFSH  
2 MMSV SQRSFWA ELNI AGLR HDNIVR VVAASTRT PEDSNS LGTI IME FGGNV TLH  
3 RSV-PC SPEAFL QEAQV - -MKKL RHEKLVQL -YAVVSEEP IYIVI EYMSKGS LLD F

S E F L E I L N L V L S N Y V I E Y G G H  
\* \* \*

1 BOV-PK - - - - - LR - R I G R F - S E P H A R F Y A A Q I V L T F E Y L H S L D L I Y R D L  
2 MMSV QVIYDATRSPEPLS CR - - KQLSLGKCLKYSLDVVNGLLFLHSQSILHLDL  
3 RSV-PC - - - - - LKGEMGKYLRLPQLVDMAAQIASGMAYVERMNYVHRDL

L R GK LSLP YAAQIV G Y HS HRDL  
\* \* \*



# Viral *src* gene products are related to the catalytic chain of mammalian cAMP-dependent protein kinase

1 BOV-PK QIEHTL NEKRI - -LQAV NFPF LVKLEFSF KDNSNLYMVM EYV PGGEMFSH  
 2 MMSV SQRSFWA ELN I AGLR HDNIVR VVAASTRT PEDSNS LG TI IME FGGNV TLH  
 3 RSV-PC SPEAFL QEAQV - -MKKL RHEKLVQL -YAVVSEEP IYIVI EYMSKGS LLD F

S E F L E I L N LV L SN Y V I E Y G G H  
 \* \* \*

1 BOV-PK - - - - - LR - R IGRF - SEPHARF YAAQIVLTF EYLHSLD LIYRDL  
 2 MMSV QVIYDATRSPEPLSCR - - KQLSLGKCLKYS LDVVNGLLFLHSQSILHLDL  
 3 RSV-PC - - - - - LK GEMGKYLR LPQLV DMAAQIAS GMAYVERMNYVHRDL

L R GK LSLP YAAQIV G Y HS HRDL  
 \* \* \*

A = Alanine  
 V = Valine  
 F = Phenylalanine  
 P = Proline  
 M = Methionine  
 I = Isoleucine  
 L = Leucine

D = Asartic Acid  
 E = Glutamic Acid  
 K = Lysine  
 R = Arginine

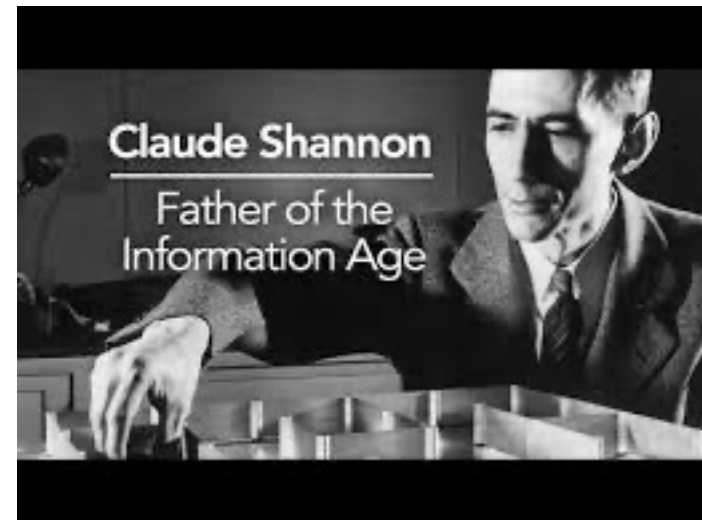
S = Serine  
 T = Threonine  
 Y = Tyrosine  
 H = Histidine  
 C = Cysteine  
 N = Asparagine  
 Q = Glutamine  
 W = Tryptophan

G = Glycine

# Information Theory

---

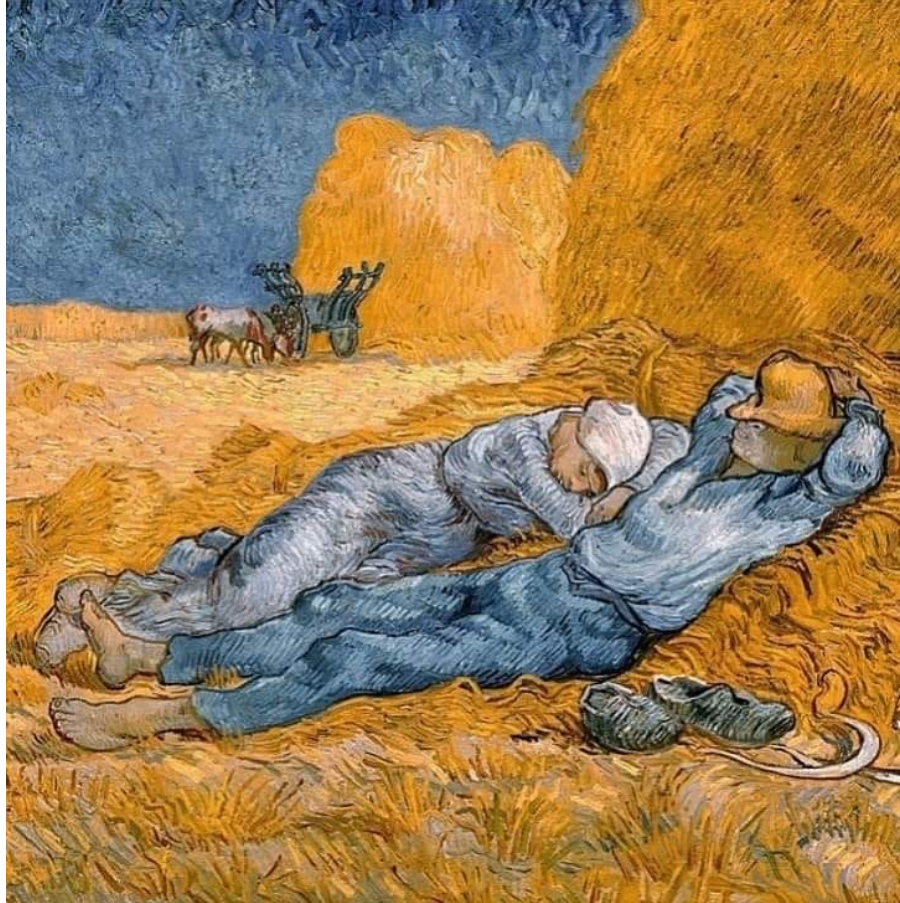
- How long an alignment should be to be statistically significant?



$$H = - \sum_{i=1}^n p_i \log p_i$$



# Chapter 2: Combinatorial Pattern Matching



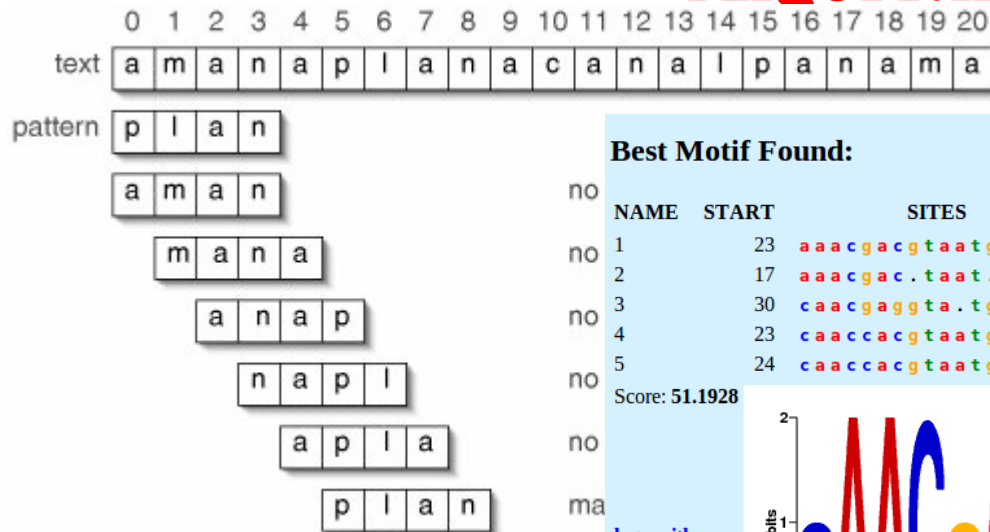
# Chapter 2

---

# Combinatorial Pattern Matching Algorithms



# Chapter 2: Combinatorial Pattern Matching Algorithms



## Best Motif Found:

no	NAME	START	SITES	END	STRAND	MARGINAL SCORE
1		23	aaacgacgtaatgctacg	6	-	22.9
2		17	aaacgac.taat.ctacg	2	-	8.45
3		30	caacgaggta.tgcaacg	14	-	14.1
4		23	caaccacgtaatgcaacg	6	-	23.6
5		24	caaccacgtaatgcatag	7	-	17.5

Score: 51.1928

[logo with ssc](#)



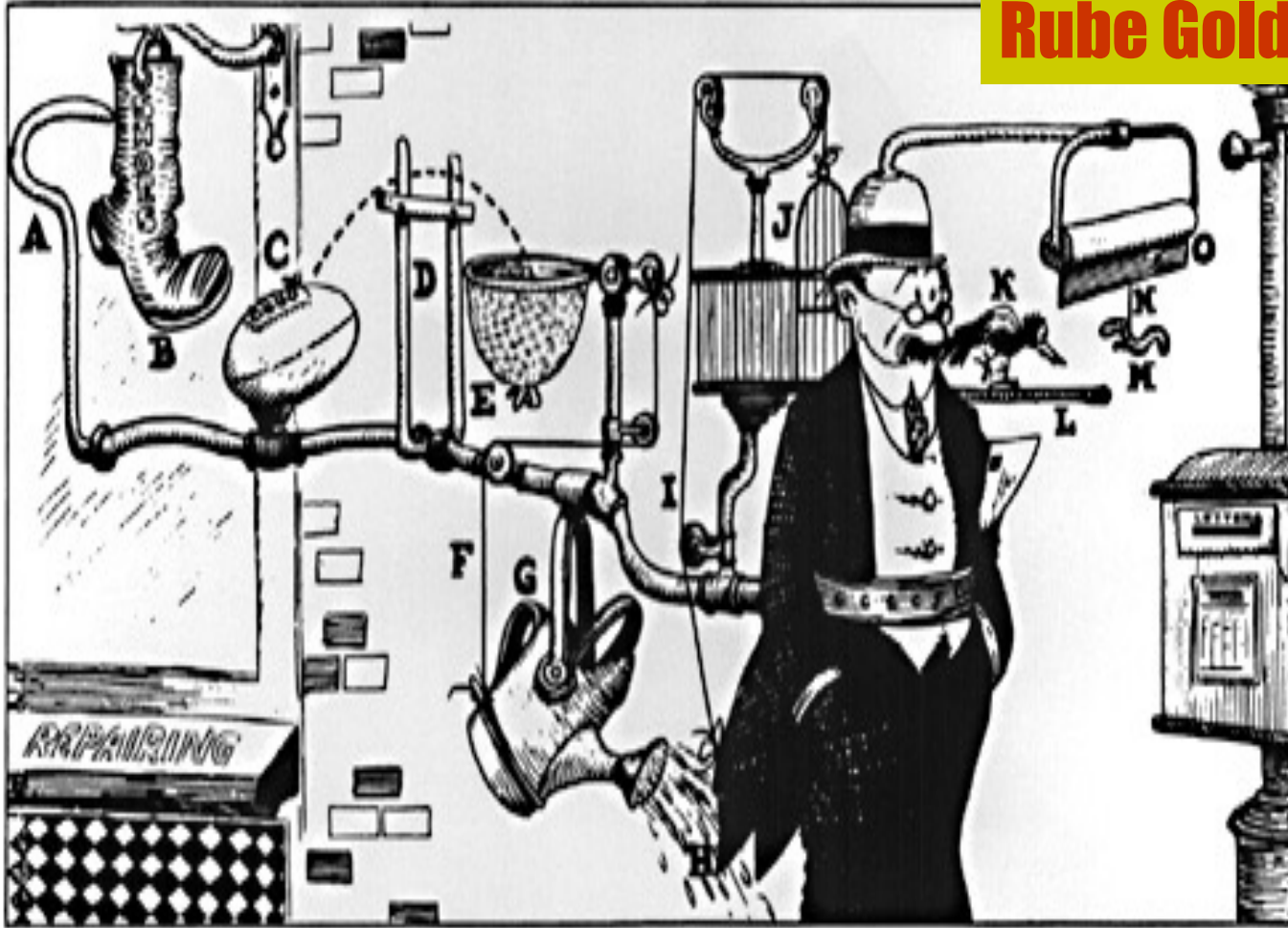
against sequence databases using [GLAM2SCAN](#).

to known motifs in motif databases using [Tomtom](#).

Regular Expression for Motif: `[ac]aac[cg]a[cg]g?taa?tg?c[at][at][ac]g`



# Rube Goldberg's Innovation



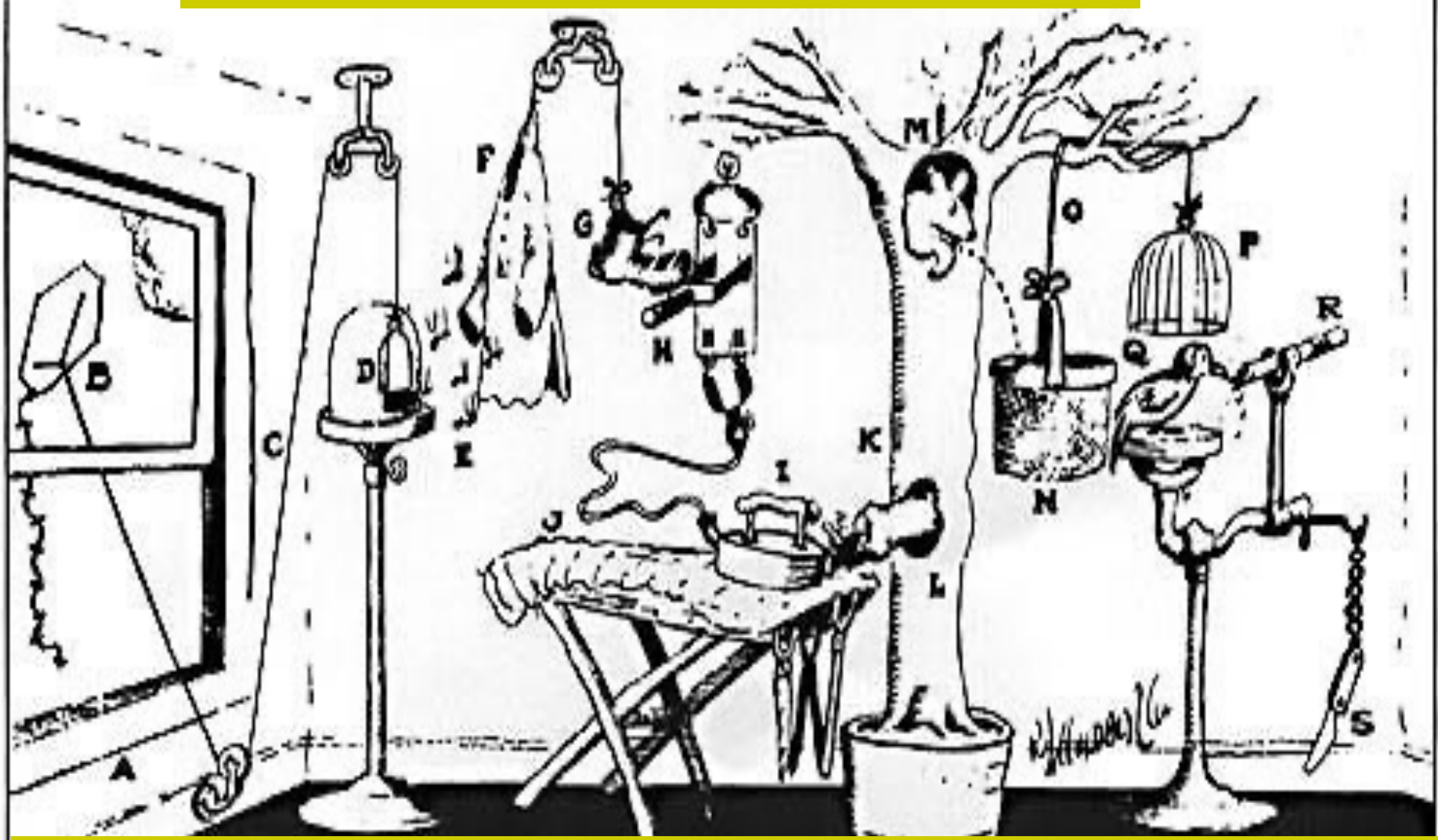
Keep You From Forgetting To Mail Your Wife's Letter RUBE GOLDBERG (tm) RGI 649

Mixed character of the problem :

continuous mathematics  
discrete mathematics

**GENOMIC  
REGULATORY  
SYSTEMS**

# Rube Goldberg's Pencil Sharpener invention



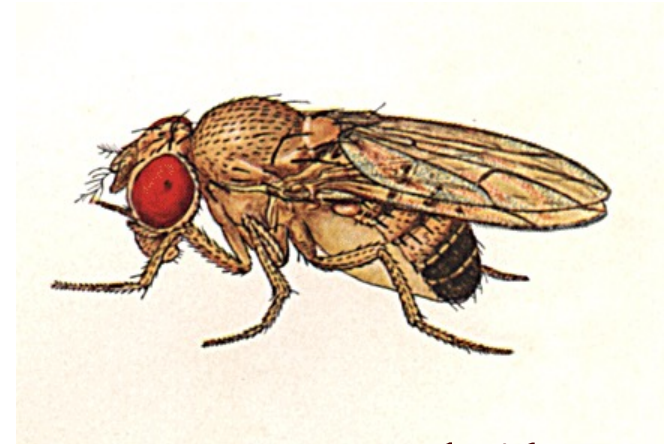
Emergency knife (S) is always handy in case opossum or the woodpecker gets sick and can't work.



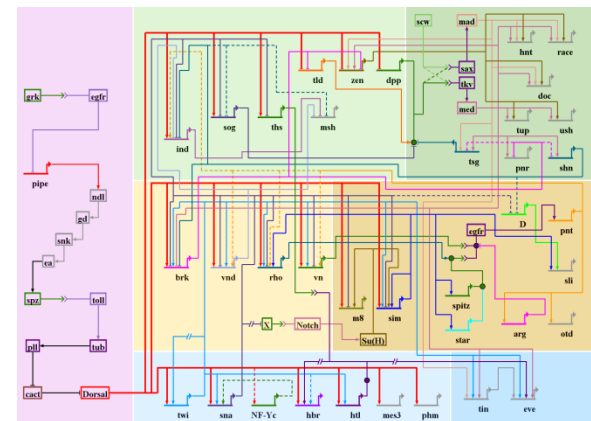
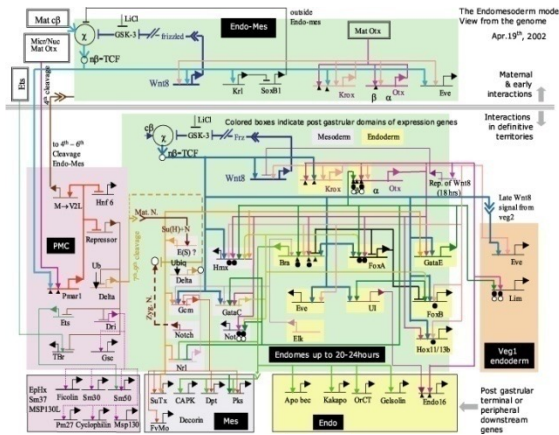
# A Tale of Two Networks



## Sea Urchin

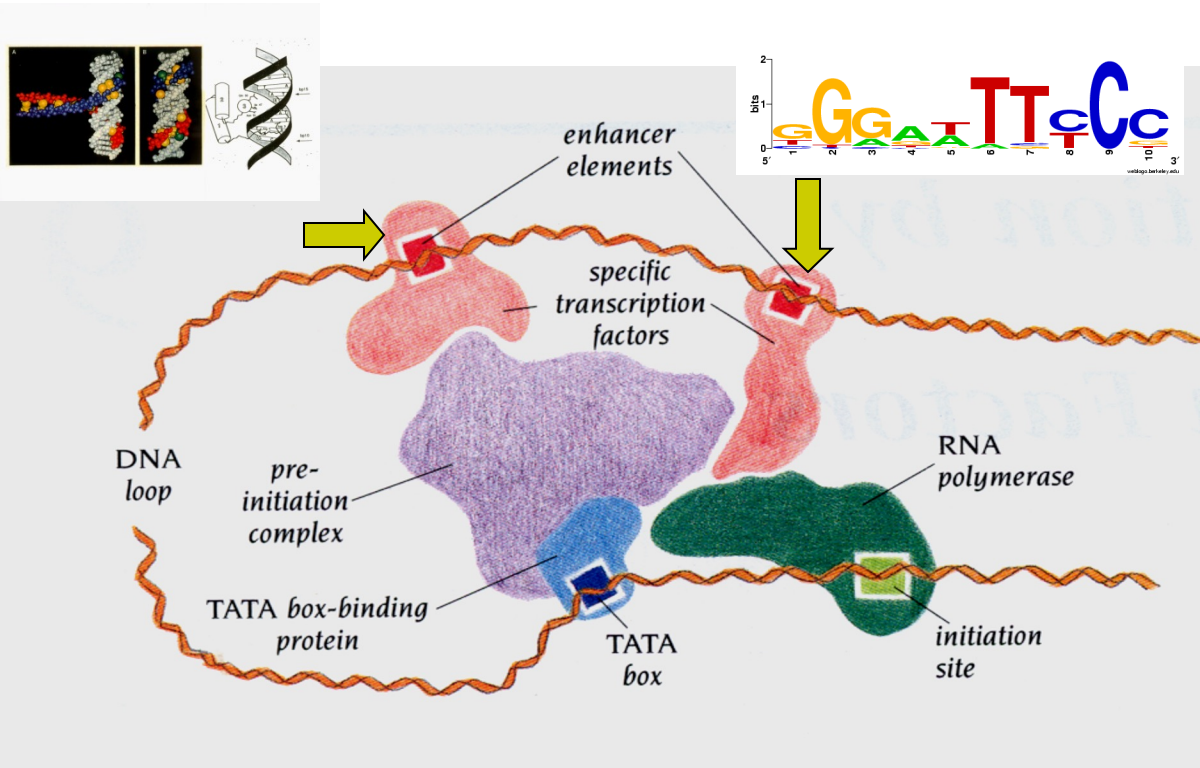


## Drosophila





# The Dogma



**Figure 9.2** Schematic model for transcriptional activation. The TATA box-binding protein, which bends the DNA upon binding to the TATA box, binds to RNA polymerase and a number of associated proteins to form the preinitiation complex. This complex interacts with different specific transcription factors that bind to promoter proximal elements and enhancer elements.





# Phylogenetic Trees (Ch. 3) ???

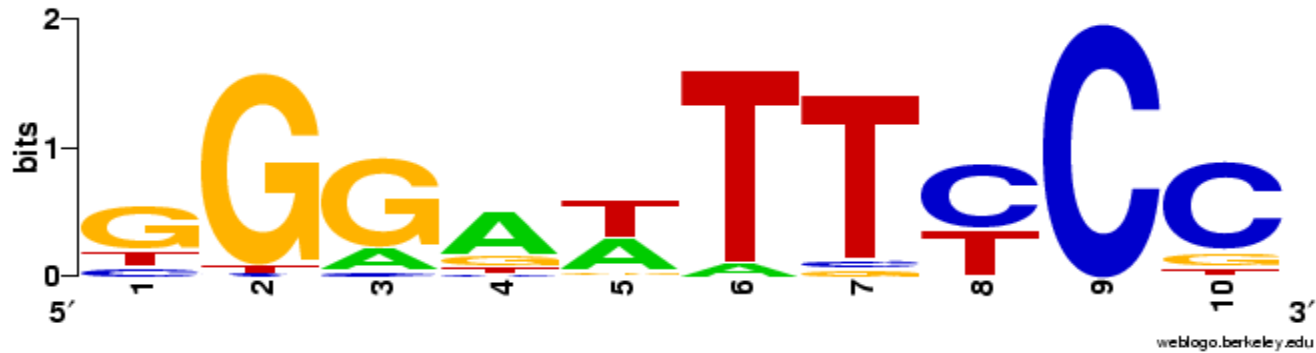
---

**Big open problem about what  
is an evolutionary model for  
regulatory regions of genes !!!**

**Phylogenetic trees are not good models for the Regulatory Genome**



# TF Binding Site Complexity









# *cis*-Regulatory Modules Complexity

A. GENE Fig. 1.2

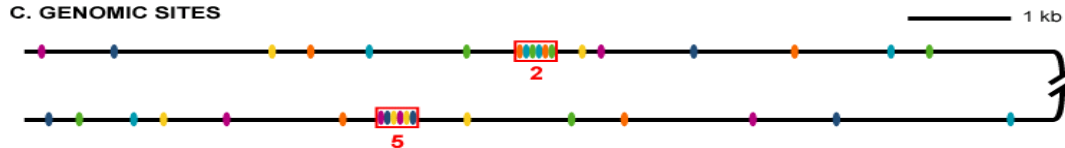
*CIS*-REGULATORY MODULES



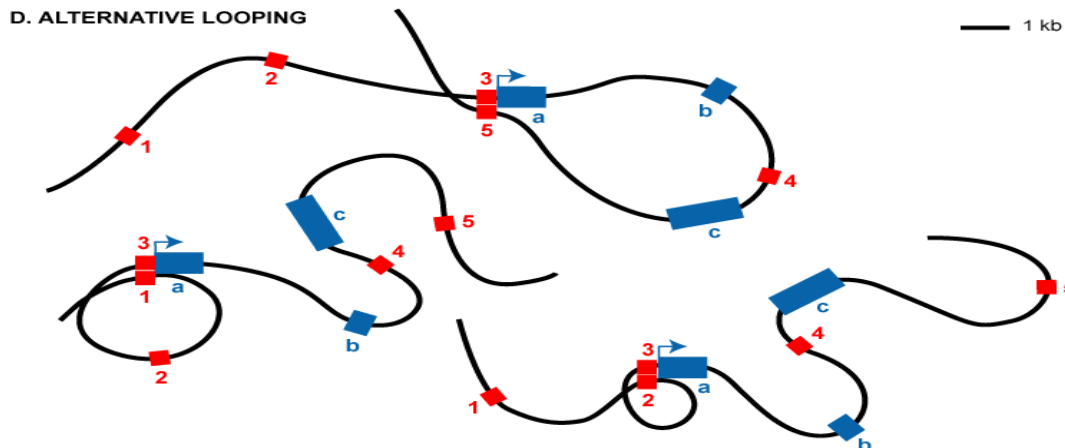
B. INPUTS



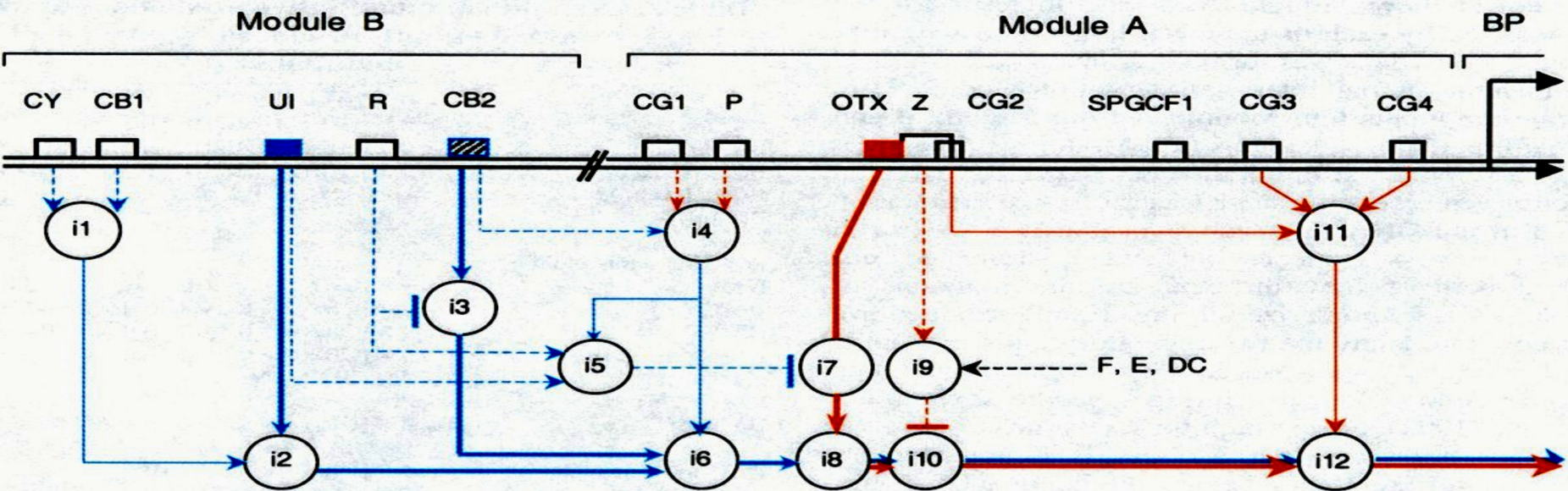
C. GENOMIC SITES



D. ALTERNATIVE LOOPING



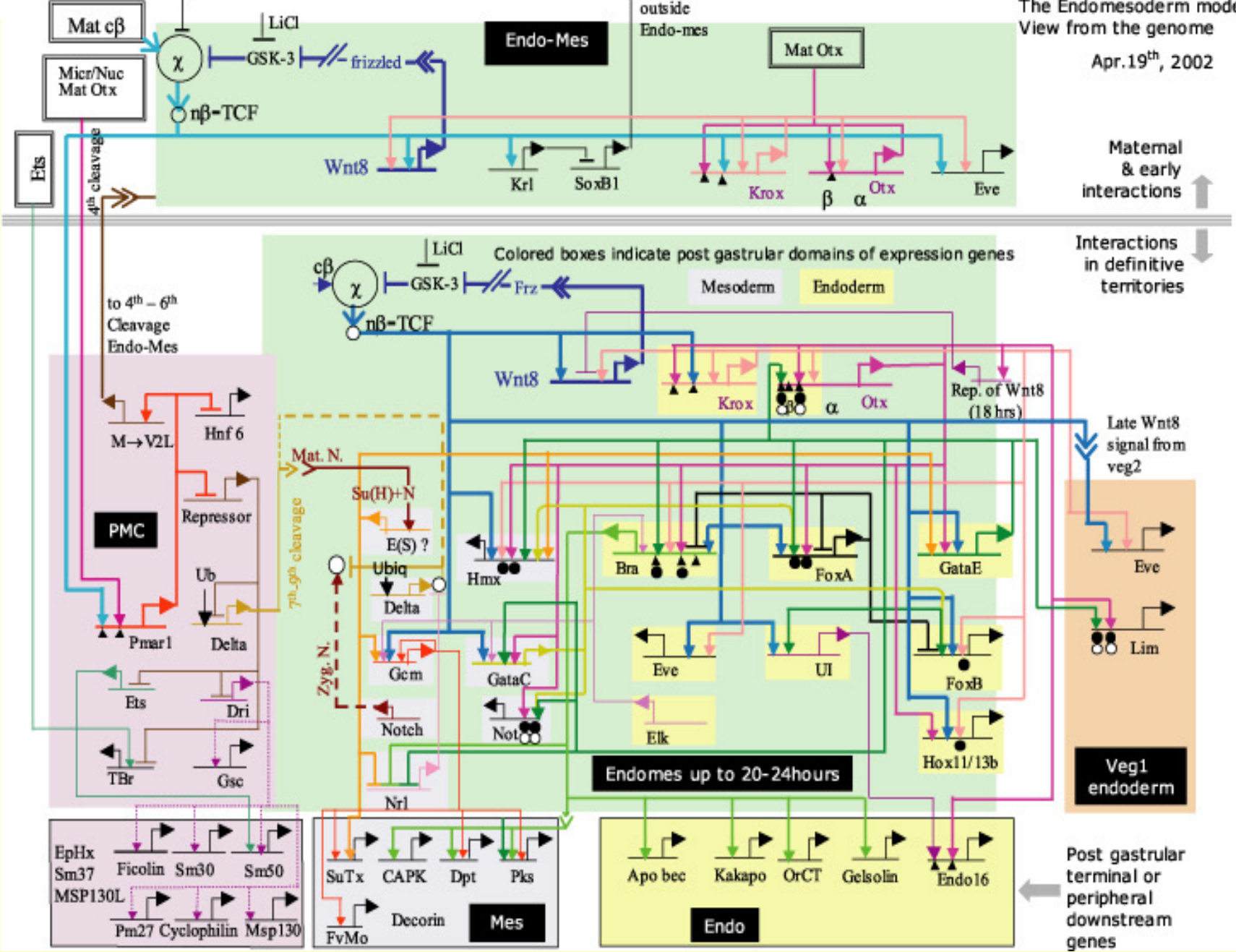
**200,000 cis-Modules**



if CY & CB1 else	$i1 = 1$ $i1 = 0.5$	if $i5 = 0$ else	$i7 = OTX(t)$ $i7 = 0$
if R else	$i2 = i1 \cdot UI(t)$	if (F or E or DC) & Z else	$i8 = i6 + i7$
if P & CG1 & CB2 else	$i3 = CB2(t)$ $i3 = k \cdot CB2(t)$ ( $1 < k < 2$ )	if $i9 = 1$ else	$i9 = 1$ $i9 = 0$
if $UI(t) > \text{threshold} \& R \& i4 \neq 0$ else	$i4 = 2$ $i4 = 0$	if (CG2 & CG3 & CG4) else	$i10 = 0$ $i10 = i8$
	$i5 = 1$ $i5 = 0$		$i11 = 2$ $i11 = 1$
	$i6 = i4 \cdot (i2 + i3)$		

The DNA program that regulates the expression of *endo16* in sea urchin

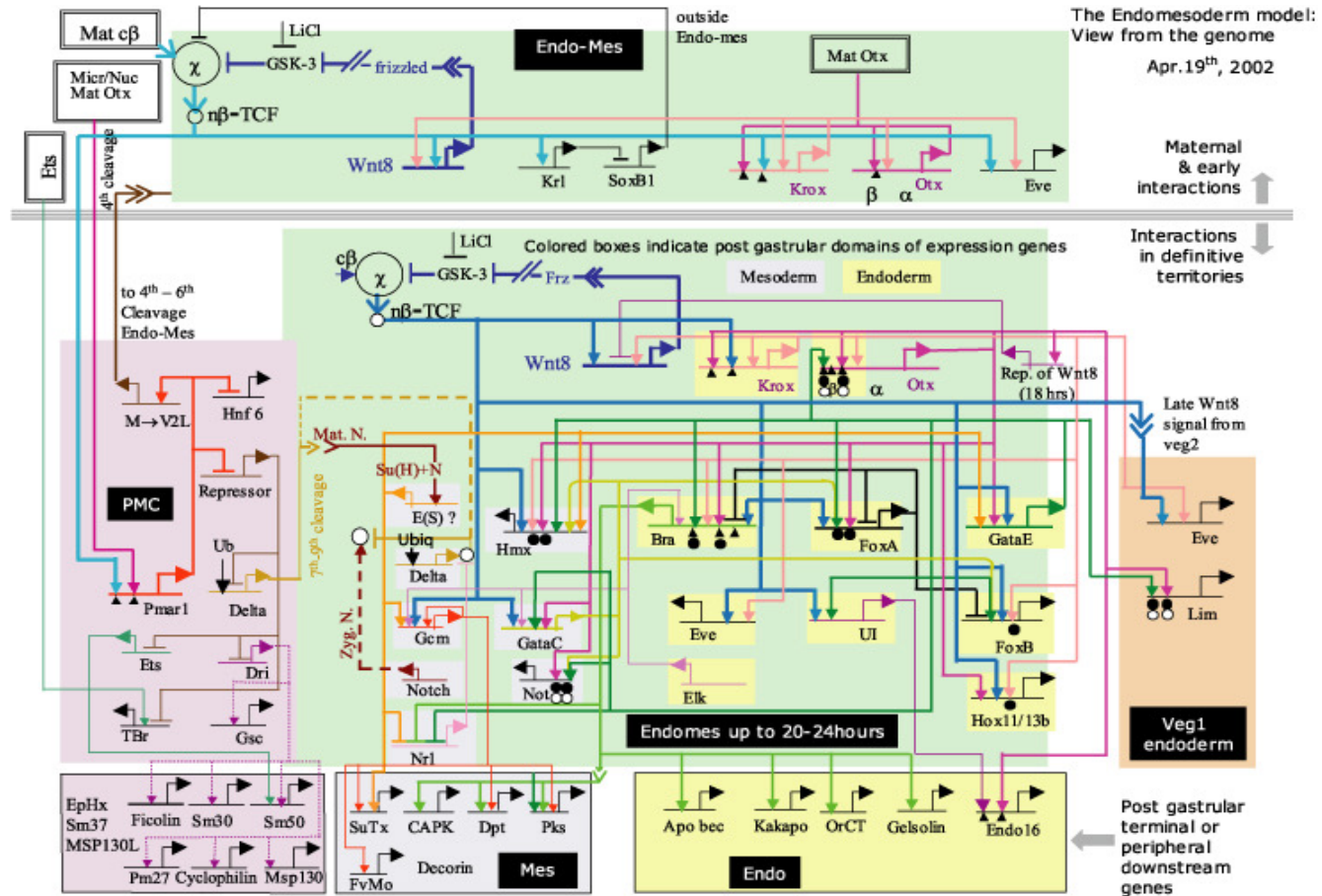
■ THE FIRST GENE



# THE FIRST NETWORK



# The View from the Genome



# A Case Study

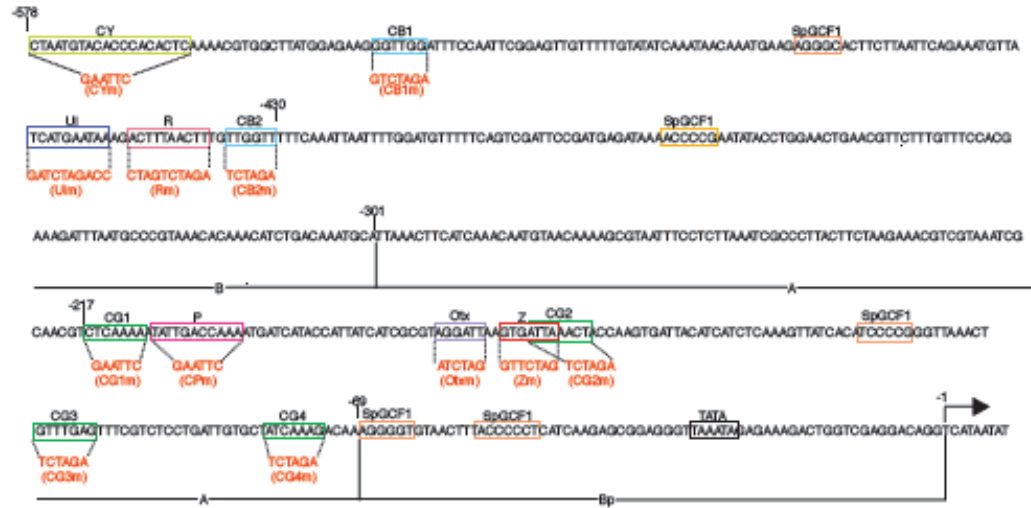


Figure 2: Quintessential diagram (from [25])

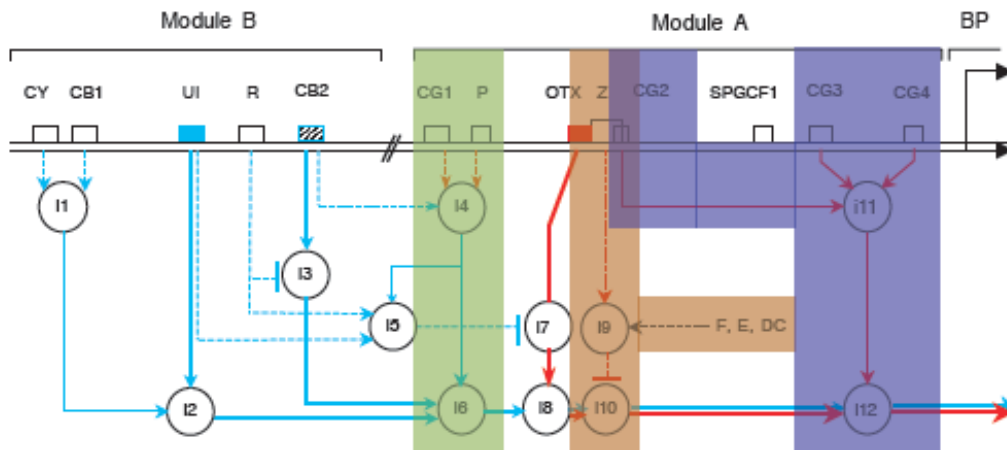


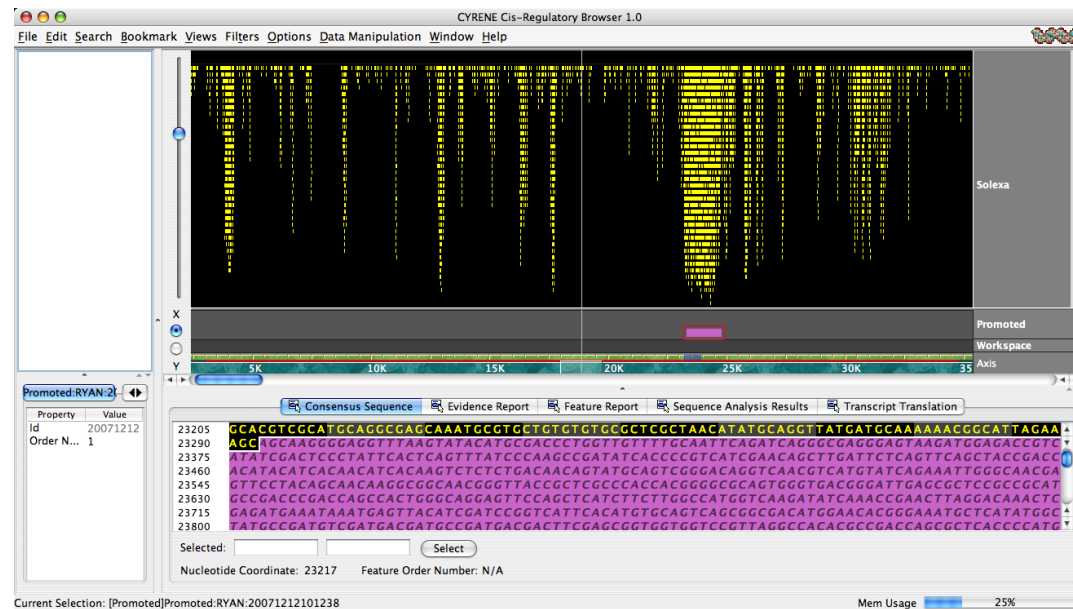
Figure 3: Computational logic model for Modules A and B of *endo16* (from [25])

# Cyrene



Ryan Tarpine

The CYRENE project seeks to address the fundamental problem of determining de novo the function of regulatory sequence by developing the cis-Lexicon, a database of known cis-regulatory modules, the cis-Browser, a next-generation regulatory genome browser, and a library of tools for assisting in the annotation pipeline. The cis-Lexicon will be a comprehensive catalog of experimentally-validated gene regulatory knowledge, designed to be a foundation and benchmark for future prediction algorithms. The cis-Browser is a high-speed integrative environment for viewing and annotating all types of genomic information. It is capable of displaying data from the cis-Lexicon, public online databases, BLAST hits, and precomputed comparative genomics analyses. To aid annotators' entry of information into the cis-Lexicon, we are developing high-throughput tools for finding relevant literature and assisting in the extraction of correct information. We suggest several algorithms to analyze the cis-regulatory data as the cis-Lexicon expands. The CYRENE project is being carried out in cooperation with Eric Davidson at the California Institute of Technology.



Cyrene Marble Probably about AD 120-150. Photo © M. Fouquet - GML.



# The cis-Browser

**Zoom Slider**

**Zoom Focus Slider**

**Annotation View**

**Selected Region**

**Loading Indicator**

**Outline View**

**Genomic Axis**

**Property Inspector View**

**Description Line**

**Selected Feature**

**Annotation Workspace Tier**

**Promoted Annotation Tier**

**Subviews**

**Popup Curation Menu**

**Memory Usage Meter**

**no sapiens:Component 4 assem**

Chr1

- GA\_x5HB7VCJ5SS (Len: 25.1)
- GA\_x54KRE8EL5L (Len: 25.3Mb)
- GA\_x54KRE8J2QS (Len: 16.79Mb)
- GA\_x5J8B7NWBLE (Len: 16.74Mb)
- GA\_x54KRE8DBDM (Len: 12.96Mb)
- GA\_x5L2HTVWYK (Len: 10.33Mb)
- GA\_x5HB7VCJ5FA (Len: 9.06Mb)
- GA\_x5L2HTU1V2W (Len: 8.4Mb)
- GA\_x5HB7VCJ5NU (Len: 8.09Mb)
- GA\_x5J8B7NYNKY (Len: 7.18Mb)
- GA\_x5J8B7P0VAE (Len: 5.98Mb)
- GA\_x54KRE8FPC6 (Len: 5.32Mb)
- GA\_x5J8B7NWR96 (Len: 4.77Mb)
- GA\_x5L2HTTM3BB (Len: 4.5Mb)
- GA\_x54KRE89K4H (Len: 4.45Mb)
- GA\_x54KRE89K4T (Len: 4.3Mb)
- GA\_x5J8B7P00BS (Len: 3.58Mb)
- GA\_x5J8B7NTU9F (Len: 3.47Mb)

2400K 2600K 2800K 3000K 3200K 3400K 3600K

hCT1951998 .WORKSPACE:34

Property	Value
CDS Links	14
Gene Accession	HCG1811273
Transcript Acc...	hCT1951998
Protein Acces...	hCP1751945
JAM Link	http://kscc191.a
Blastx Nraa Pr...	http://genedisc
Id	CELERA:110005
Order Number	0
Comments	1
Curation Flags	2
Feature Type	Transcript
Algorithm: Data...	Curation
Parent Feature Id	WORKSPACE:33
Axis Name	GA_x5HB7VCJ5
Axis Id	CELERA:1950001
Axis Begin	2804129
Axis End	3248807

Frame +1 0 M E N L I R G G R N P P Q Y Q R S

Frame +1 50 A G G A T G G A G A A C C T A A T A A G G G A A G G A A T C C C C C G C A A T A C C A G A G A G A

Frame +1 100 T C C T T G T A A A G A G G T T C G T C A G C A C T T C G G A A G A G G C C T G A A G A G A G A G A G

Frame +1 150 G C T A G C T G A G A G A G A C C A C C A A E R K I I G I G G C C C C C T G G T

Frame +1 200 T G C A G T C G T G A G A G A C C A C C A A E R K I I G I G G C C C C C T G G T

Frame +1 250 P C R C F R G E E E E T C C G A P

Frame +1 300 C C A T G C C G A T T T T C C G A G G T G A A G R K I I G I G G C C C C C T G G T

Frame +1 350 S H C S L Q G V T C G A A

Frame +1 400 C T C C C A C T G C A G T T T C A G C A G G T G

Frame +1 450 T L E E L Y L R L D T C G A A

Frame +1 500 G A A C A T T A G A G G A C T T T A T C T A G A

Frame +1 550 K Q F N C Q A A L

Frame +1 600 A A G C A A T T G T T C A A C T G T C A A G C T C

Frame +1 650 D L S N L P T

Translation Frames Selected: +1 ORF

Nucleotide Coordinate: 126 Amino acid Coordinate: 41

Selected Range: 122 : 122 Selected Amino Acid Length: 0

Current Selection: [Curation] hCT1951998 .WORKSPACE:34

Mem Usage

**Curation Options:**

- Set Start Codon
- Set Translation Start
- Set Longest Open Reading Frame
- Set Longest ATG to Stop
- Delete Start Codon
- Delete Start Codon From Database
- Set Stop Codon to Calculate ORF upstream

# Transcript Curation

**Search Known Features**

Search:  
 Component 4 assembly from 2001-03-08  
 1 loaded Genome Version(s)  
 5 available Genome Version(s)

Type: Gene (Celera Accession) [v]  
Find: hCG1811273  
[Search]

Status: Search Complete

Results:  
Gene hCG1811273

Open New Browser [Bookmark] [Navigate] [Stop] [Close]

Order Number	0
Comments	1
Curation Flags	2
Feature Type	Transcript
Algorithm:Dataset	Promoted
Parent Feature Id	CELERA:1100059
Axis Name	GA_x5HB7VCJ5S
Axis Id	CELERA:1950000
Axis Begin	2804129
Axis End	3248807

Current Selection: [Promoted] hCT1951988 :CELERA:11000595808682

Mem Usage [Progress Bar]

Tracks (right side): S4:human\_dbEST, Bn:mouse, Bn:human\_dbEST, Bn:dog, ConsSegment, Splice Donor, Splice Acceptor, FgenesH, Otto, Promoted, Workspace, Axis, Workspace (rev), Promoted (rev), Otto (rev), SNP (rev), Genscan (rev), FgenesH (rev), Splice Acceptor (rev), start:1 (rev), start:2 (rev), start:3 (rev), stop:1 (rev), stop:2 (rev), stop:3 (rev), Splice Donor (rev), Bn:dog (rev), Bn:human\_dbEST (rev), Bn:mouse (rev), S4:human\_dbEST (rev), Bn:CMGI (rev), Bx:nraa (rev), Bn:CHGI (rev), GRILL (rev)



# Sequence Comparison

The screenshot displays a bioinformatics software interface with the following components:

- Top Panel:** A genomic map showing various tracks. The left sidebar lists sequences such as GA\_x5HB7VCJ5SS (Len: 25.3Mb) and GA\_x54KRE8EL5L (Len: 16.79Mb). The right sidebar shows tracks for GRAIL, Bn:CHGI, Bx:nraa, Bn:CMGI, S4:human\_dbEST, Bn:mouse, Bn:human\_dbEST, Bn:dog, Otto, Promoted, Workspace, Workspace (rev), Otto (rev), and SNP (rev).
- Bottom Panel:** A detailed view of a consensus sequence and its alignments. The top part shows a table of properties for the selected sequence (Bx:nraa:CELERA:50000090728954). The middle part shows a sequence alignment viewer with a consensus sequence and several alignments. The bottom part shows a table of properties for the selected sequence (Bx:nraa:CELERA:50000090728954).

**Table of Properties (Left Panel):**

Property	Value
Id	CELERA:50000090728954
Aliases (Numb...)	0
Order Number	1
Comments	0
Feature Type	High Scoring Pair
Algorithm:Data...	Bx:nraa
Parent Feature...	CELERA:50000090728954
Axis Name	GA_x5HB7VCJ5
Axis Id	CELERA:1950001
Axis Begin	5261006
Axis End	5261183
Entity Length	177
Entity Orientation	Forward
Is Child	true
Is Composite	false
Relative Asse...	985037830
Display Priority	Low

**Table of Properties (Bottom Panel):**

Property	Value
Id	CELERA:50000090728954
Aliases (Numb...)	0
Order Number	1
Comments	0
Feature Type	High Scoring Pair
Algorithm:Data...	Bx:nraa
Parent Feature...	CELERA:50000090728954
Axis Name	GA_x5HB7VCJ5
Axis Id	CELERA:1950001
Axis Begin	5261006
Axis End	5261183
Entity Length	177
Entity Orientation	Forward
Is Child	true
Is Composite	false
Relative Asse...	985037830
Display Priority	Low

**Sequence Alignment Viewer (Bottom Panel):**

Consensus Sequence: K K E L T Q I K Q K

Alignment 1: K K E L T Q I K Q K (+3)

Alignment 2: K K E L T Q I K Q K (+2)

Alignment 3: S G Q R G S S K L K G D D L Q A I K K (+1)

Alignment 4: S G Q V G S S E L K G D D L Q A I R R (+1)

DNA +

Query Sequence: GTGGGCA GTGGGATCTTCCGAATTGAAAGGATGACCTTCA GGCCATAA GAA GGAATT GACCA GAT AAAA CAAAA

Axis: 5260960 5260970 5260980 5260990 5261000 5261010 5261020 5261030

DNA -

-1

-2

-3

CYRENE Cis-Regulatory Browser 1.0

File Edit Search Bookmark Views Filters Options Data Manipulator Window Help

Strongylocentrotus p  
Unknown Chromos

endo16

Promoted

Workspace

Axis

10K 15K 20K

373279 :ENDO:12345

Property	Value
Gene ...	373279
Id	123456789
Descri...	Polyfunctio
Alias ...	endo16

Consensus Sequence Evidence Report

```

11798 TAAATAGAGAAAGACTGGTCGAGGACAGGTCATA
11832 ATATTGCTAATTTTGGAGACGATGAGGAGGTTAA
11866 ATATTTTGCTGTTTCGCGGTTTGGCCGTGGCGCG
11900 GTCAATGCCACAGGTAAGAAATATAATAATTT
11934 ACAAATTAGTTTTAAGACGCCTCCTTCTTCTTC
11968 TTGCTTTTCAACTGTTAAATACATGCTTTTGTG
  
```

Selected: 11853 11855 Select 66

Nucleotide Coordinate: 11902 Feature Order Number: N/A

Current Selection: [Promoted] 373279 :ENDO:123456789 Mem Usage 6%

# Inter-species comparison

The screenshot displays the CYRENE Cis-Regulatory Browser 1.0 interface. The main window shows a genomic axis with three species' cis-regulatory regions (CFTR) aligned: Homo sapiens (red), Baboon (purple), and Chicken (green). The axis is labeled from 500K to 750K. Below the axis, a detailed view of a specific alignment is shown, with the sequence: TTTCTCAGGAATCACTGACATAGGAGAAGTTTCCCAATTTCTGACCGAGGGAATCATCATGAAAGATTTTAGTCATCC. The alignment is shown with a score of 79.7% identity and 100.0% coverage. The interface includes a menu bar (File, Edit, Search, Bookmark, Views, Filters, Options, Data Manipulation, Window, Help), a toolbar, and a status bar at the bottom showing 'Current Selection: [At:chicken\_CFTR\_anno]At:chicken\_CFTR\_anno:CELERA:30000103836821' and 'Mem Usage 21%'.

**Species and Regions:**

- Homo sapiens cfr region: Humana
- Baboon cfr region: Baboon CFTR
- Chicken cfr region: Chicken CFTR

**Genomic Axis (Len: 1.88M):**

- Unknown Chromosome
- Genomic Axis (Len: 1.88M)

**Axis Labels:**

- At:chicken\_CFTR\_a
- chicken\_CFTR\_axis
- At:baboon\_CFTR\_a
- baboon\_CFTR\_axis
- Curated:human\_CFTR
- Axis
- Workspace

**Entity Properties:**

Property	Value
Id	3000010383682
Subject Left	281798
Subject Right	281980
Alignment Len...	182
Aliases (Numb...	0
Feature Type	High Scoring Pair
Algorithm:Data...	At:chicken_CFTR_...
Axis Id	HUMT1:0
Axis Begin	608309
Axis End	608491
Entity Length	182
Entity Orientati...	forward
Relative Assem...	1
Display Priority	high
Assembly Vers...	1
Exon Name	CHKT1IMET-17
Mapping Status	URU
Percent Covera...	100.0
Percent Id	79.7

**Sequence Analysis Results:**

Sequence	Score
gAAATACGAGTTCAGAAAGTTTCCTGAGGAAATATCATGAAAGATTTATTCATCC	+3
gAAATACGAGTTCAGAAAGTTTCCTGAGGAAATATCATGAAAGATTTATTCATCC	+2
gAAATGACTGACATAGGAGAAGTTTCCTGAGGAAATATCATGAAAGATTTAGTCATCC	+1
gAAATGACTGACATAGGAGAAGTTTCCTGAGGAAATATCATGAAAGATTTAGTCATCC	+1

**Query Sequence:** TTTCTCAGGAATCACTGACATAGGAGAAGTTTCCCAATTTCTGACCGAGGGAATCATCATGAAAGATTTTAGTCATCC

**Axis Labels (Detailed View):**

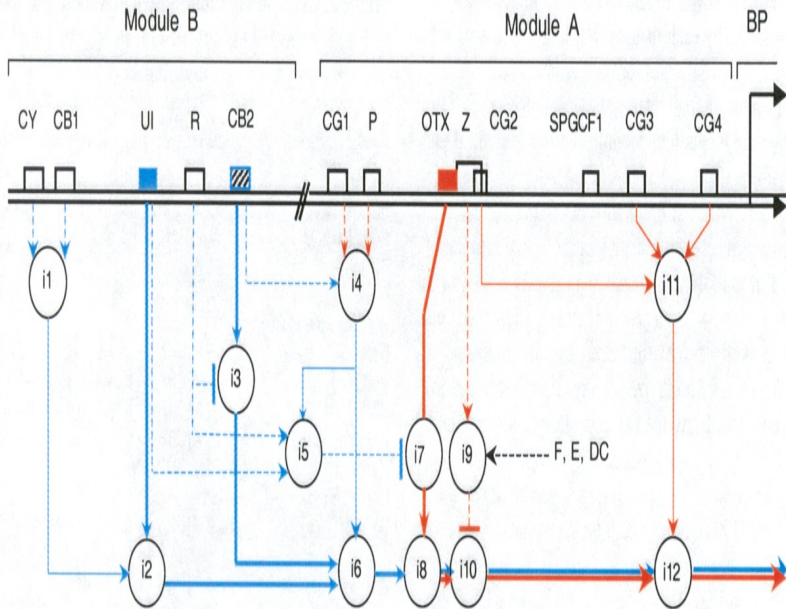
- DNA +
- Query Sequence
- Axis
- DNA -
- 1
- 2

**Current Selection:** [At:chicken\_CFTR\_anno]At:chicken\_CFTR\_anno:CELERA:30000103836821

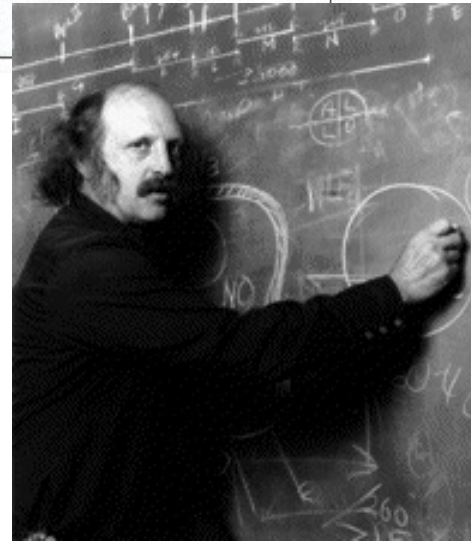
**Mem Usage:** 21%

# One gene, 30 years of study, 300 docs and postdocs

## A Proposal for Nobel Prize



if CY & CB1	$i1 = 1$	if $i5 = 0$	$i7 = OTX(t)$
else	$i1 = 0.5$	else	$i7 = 0$
			$i8 = i6 + i7$
	$i2 = i1 \cdot UI(t)$		
		if (F or E or DC) & Z	$i9 = 1$
if R	$i3 = CB2(t)$	else	$i9 = 0$
else	$i3 = k \cdot CB2(t)$ ( $1 < k < 2$ )	if $i9 = 1$	$i10 = 0$
		else	$i10 = i8$
if P & CG1 & CB2	$i4 = 2$	if (CG2 & CG3 & CG4)	$i11 = 2$
else	$i4 = 0$	else	$i11 = 1$
if $UI(t) > \text{threshold} \& R \& i4 \neq 0$	$i5 = 1$		
else	$i5 = 0$		
			$i12 = i11 \cdot i10$
	$i6 = i4 \cdot (i2 + i3)$		



“Programs built into the DNA of every animal.”

Eric H. Davidson

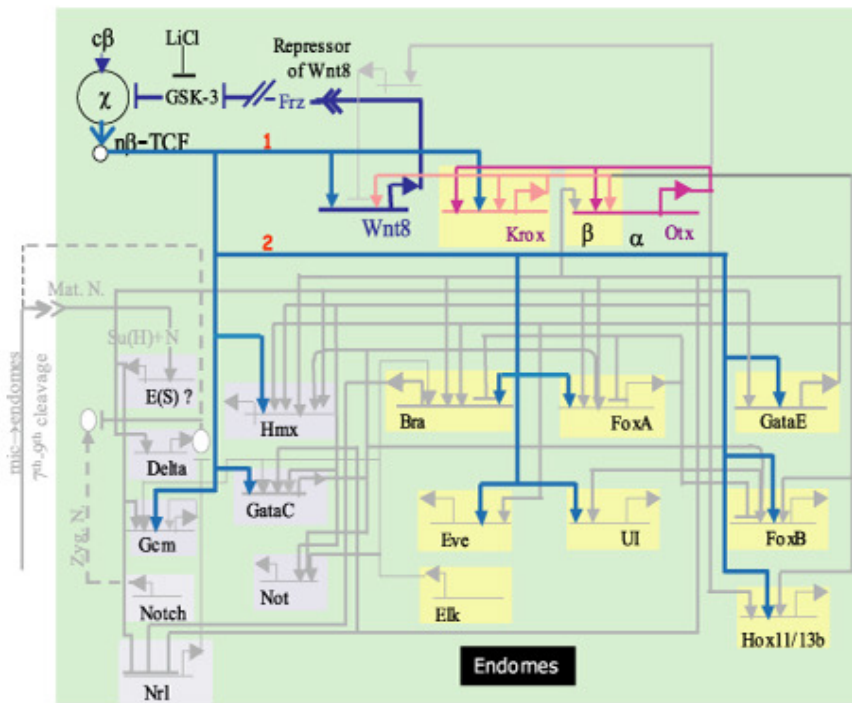
# Genomic Regulatory Systems



# The View from the Nucleus

View from the nucleus: Endomesoderm nuclei to hatching blastula stage; the Wnt8/Tcf signalling loop and its genes.

Apr. 19<sup>th</sup>, 2002



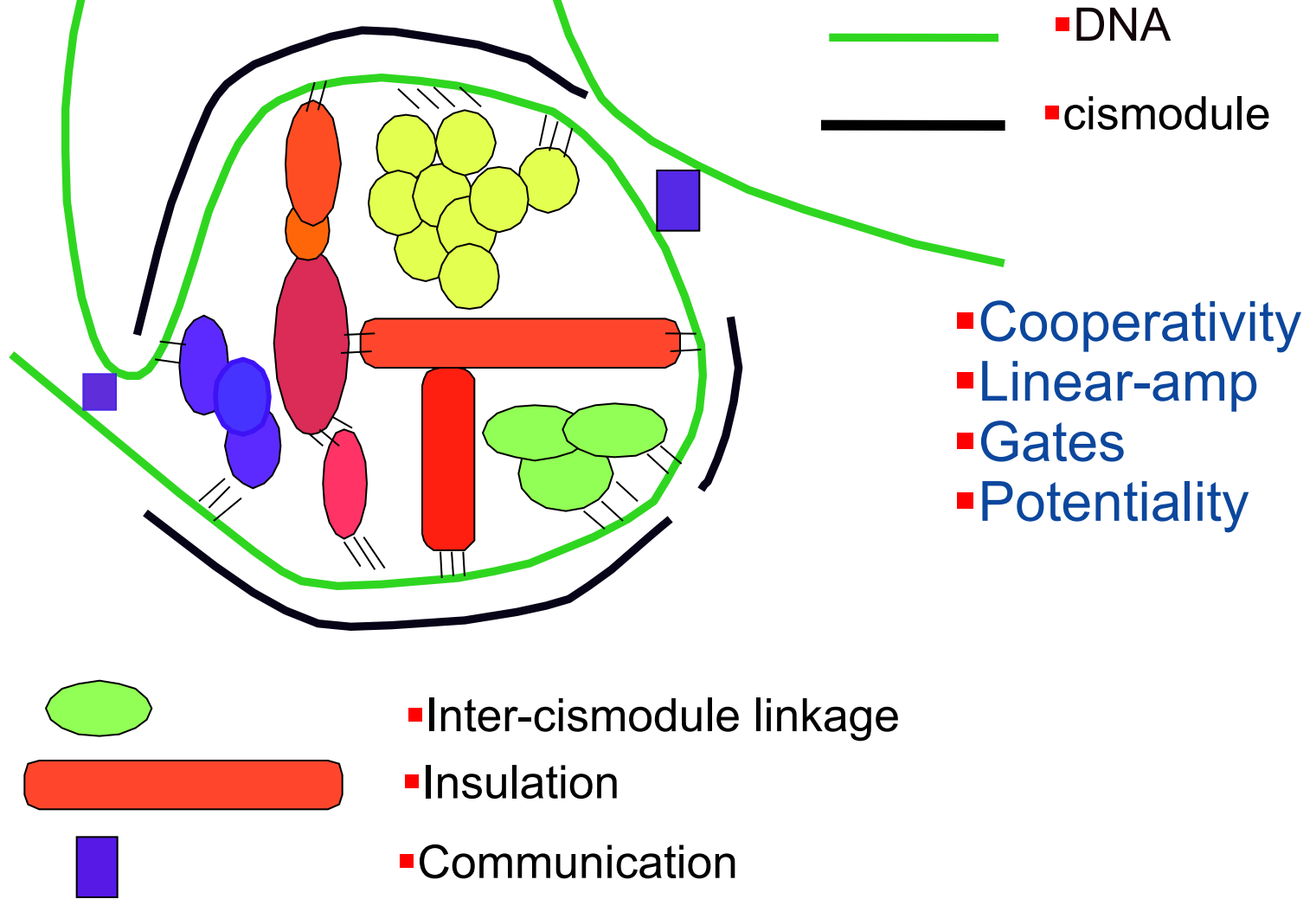
Notes:

1.  $\beta$ -catenin/Tcf input now produced by a zygotic signaling loop driven by Wnt8 expression in endomesoderm cells.
2.  $\beta$ -catenin/Tcf input required for expression of many regulatory genes that become active in the  $veg_2$  endomesodermal territory during early-mid blastula stage.

	Post gastrula mesoderm only		Post gastrula endoderm only		Endomesoderm upto 20-24 hours
--	-----------------------------	--	-----------------------------	--	-------------------------------

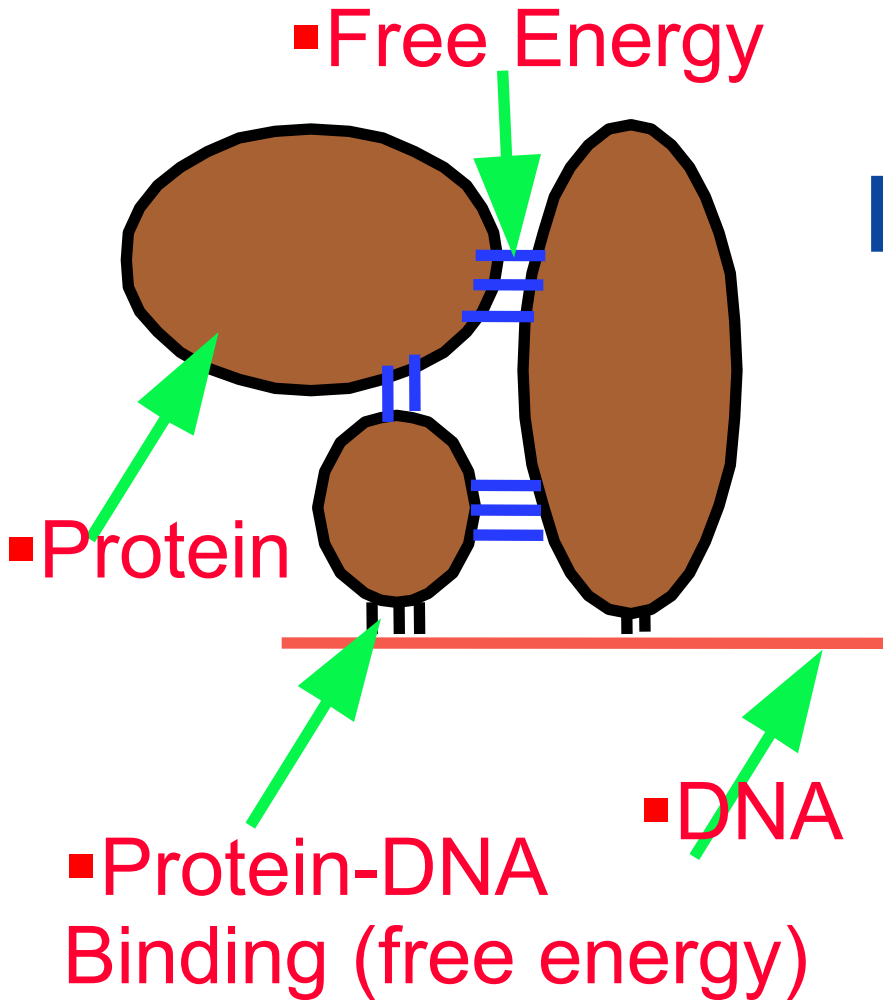


# Building Protein-DNA Assemblies





# The Building Blocks



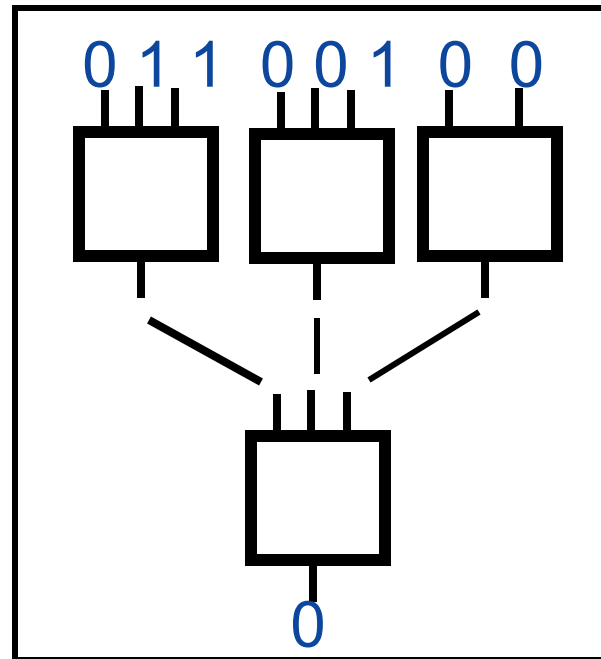
**Free energy is the “GLUE”**



---

# Information Processing

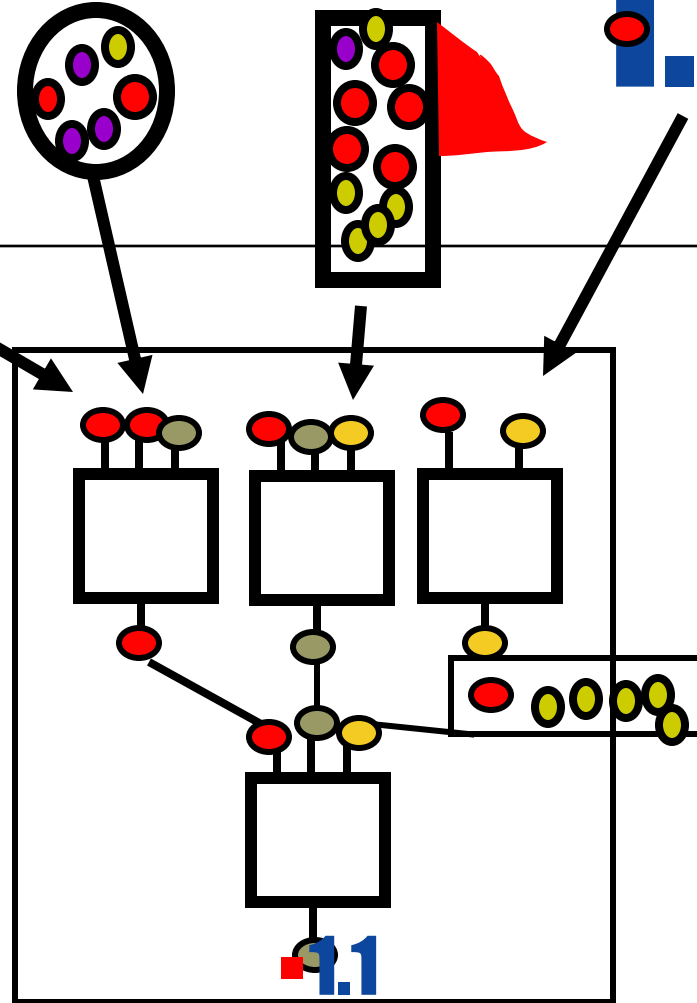
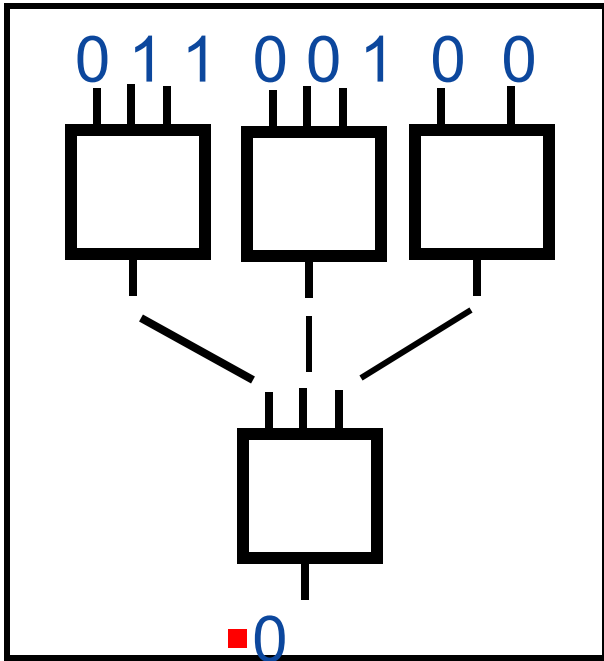




- **Boolean Circuit**
- **Synchronous** input and output
- **Completely** defined gates

0.5

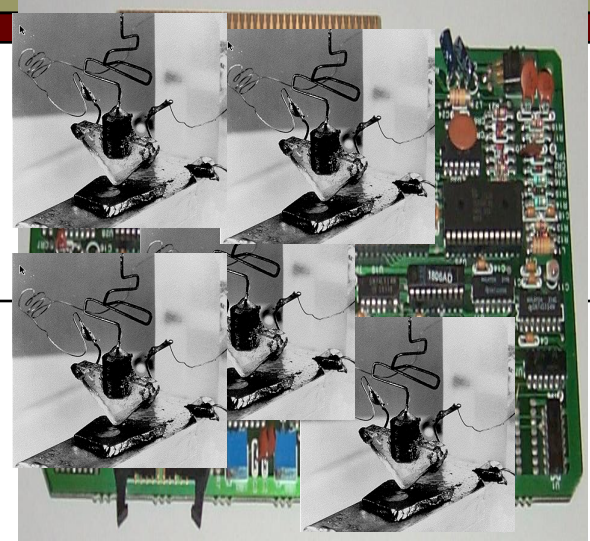
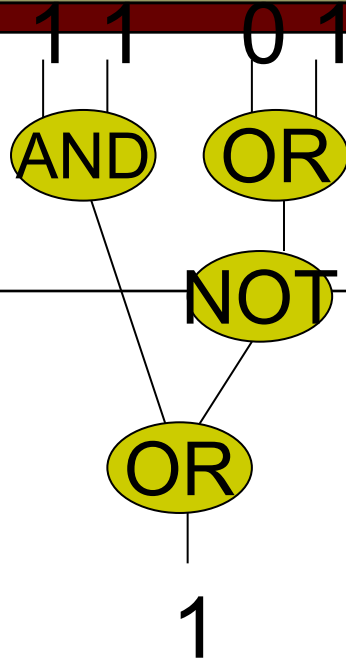
1.4



- **Boolean Circuit**
- **Synchronous** input and output
- **Completely** defined gates

- **Boolinear Circuit**
- **Asynchronous** input and output
- **Incompletely** defined gates

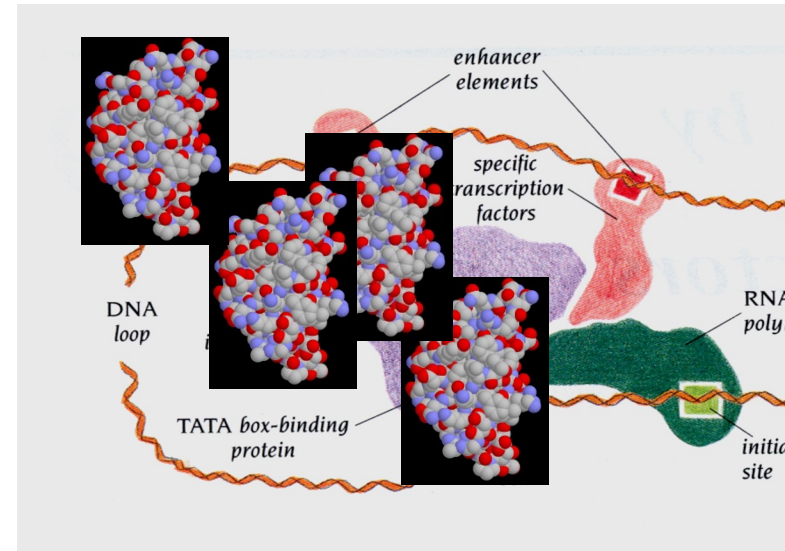
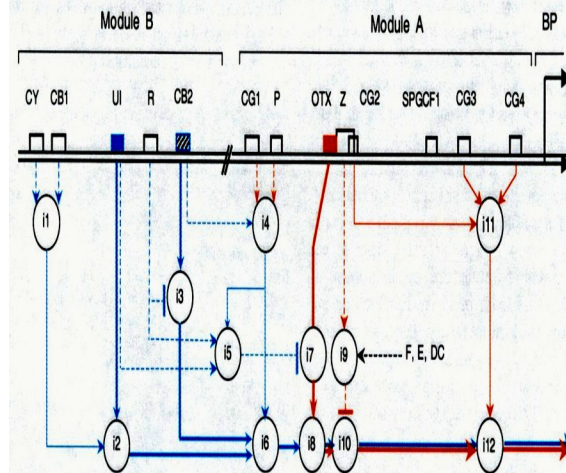
IF (x1 = 1  
 AND  
 x2 = 1)  
 THEN  
 .....



GTAGGATTAAG  
 .....

CATCCTAATTC

GTATCTAGAAG  
 .....



# Chapter 3: Phylogenetic Trees

---



Image, courtesy of Vincent van Gogh Museum

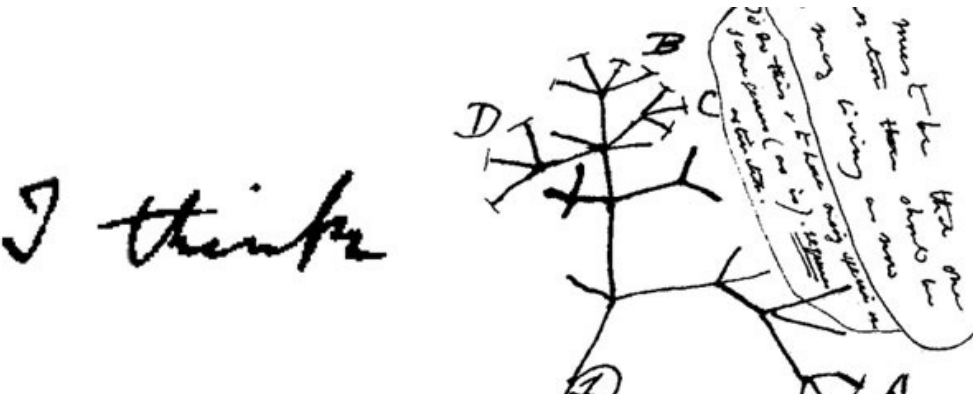
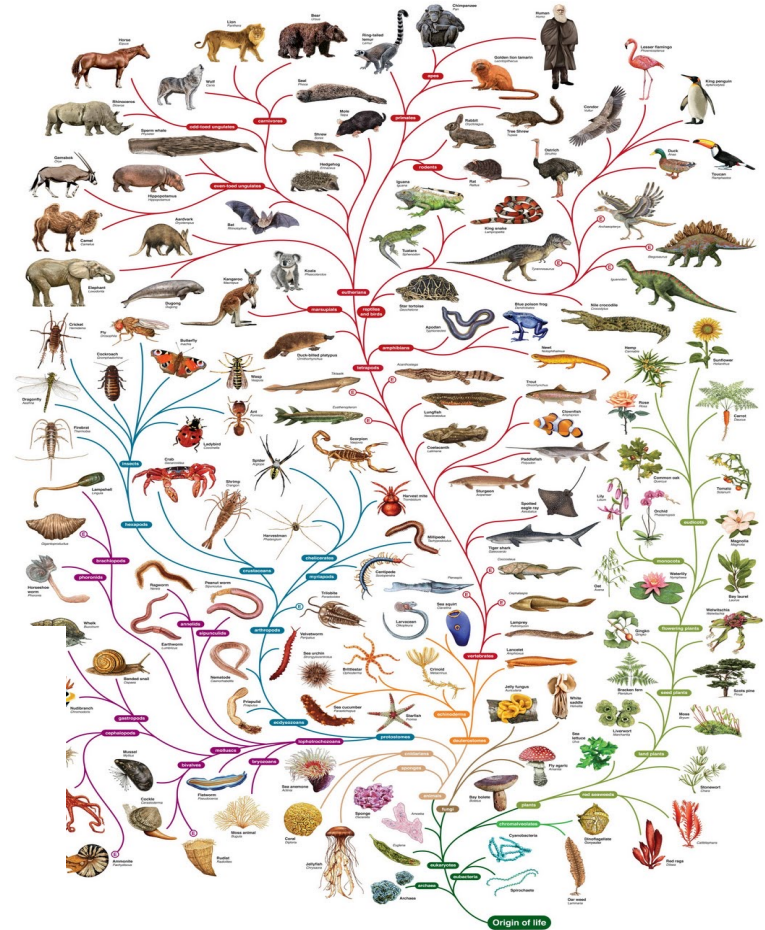


# Chapter 3

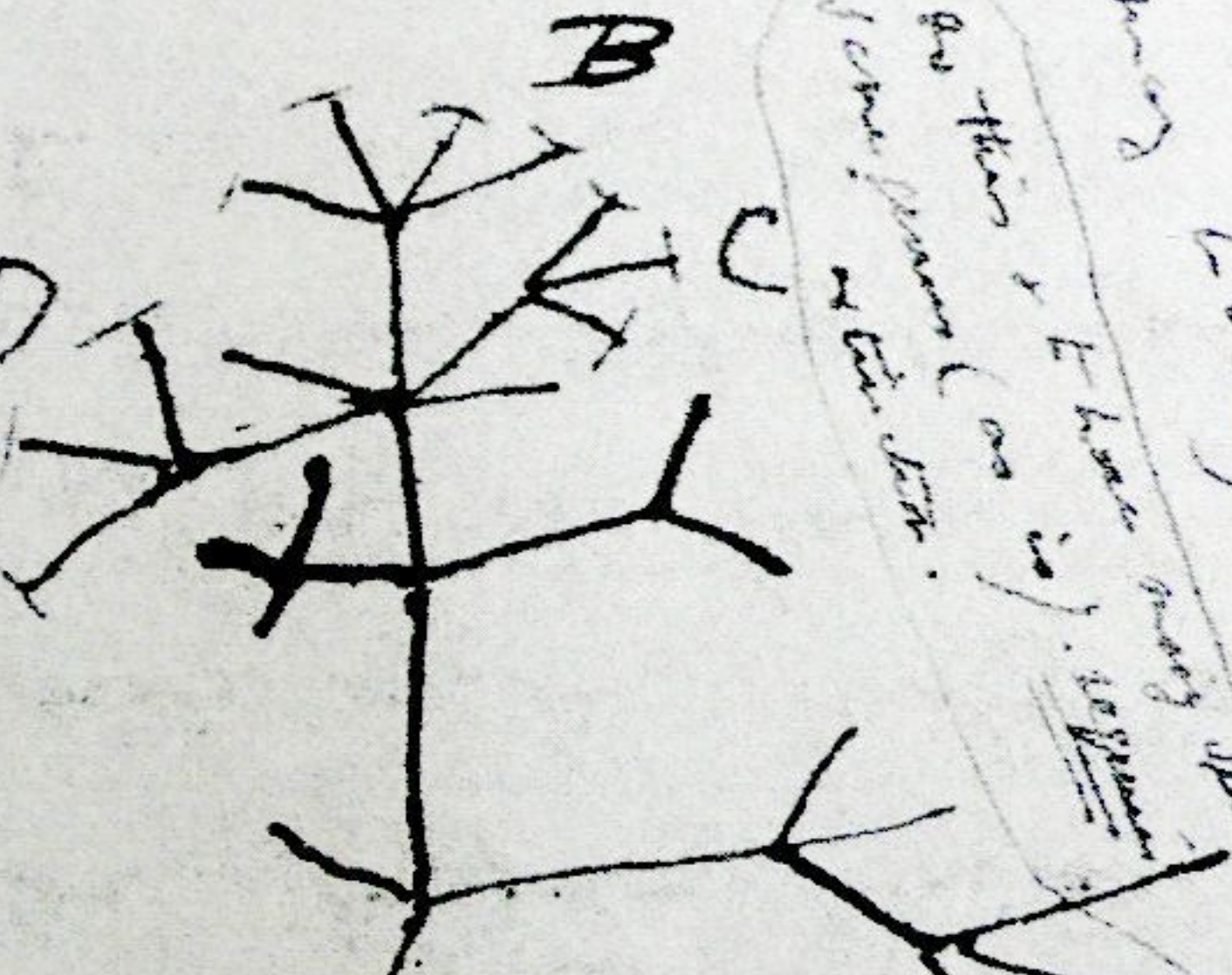
---

# Phylogenetic Trees Algorithms

# Chapter 3: Phylogenetic Trees Algorithms



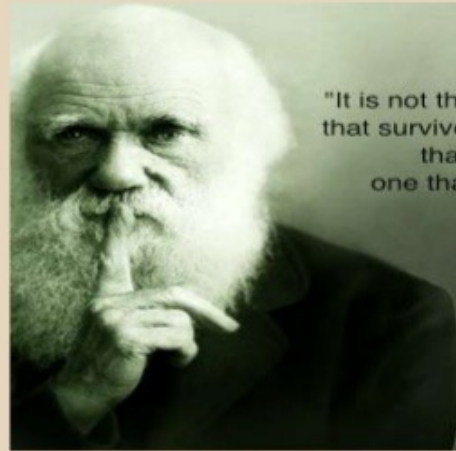
thanks



Do the things & E have many species  
 as the same species (as in) the same  
 as the same.

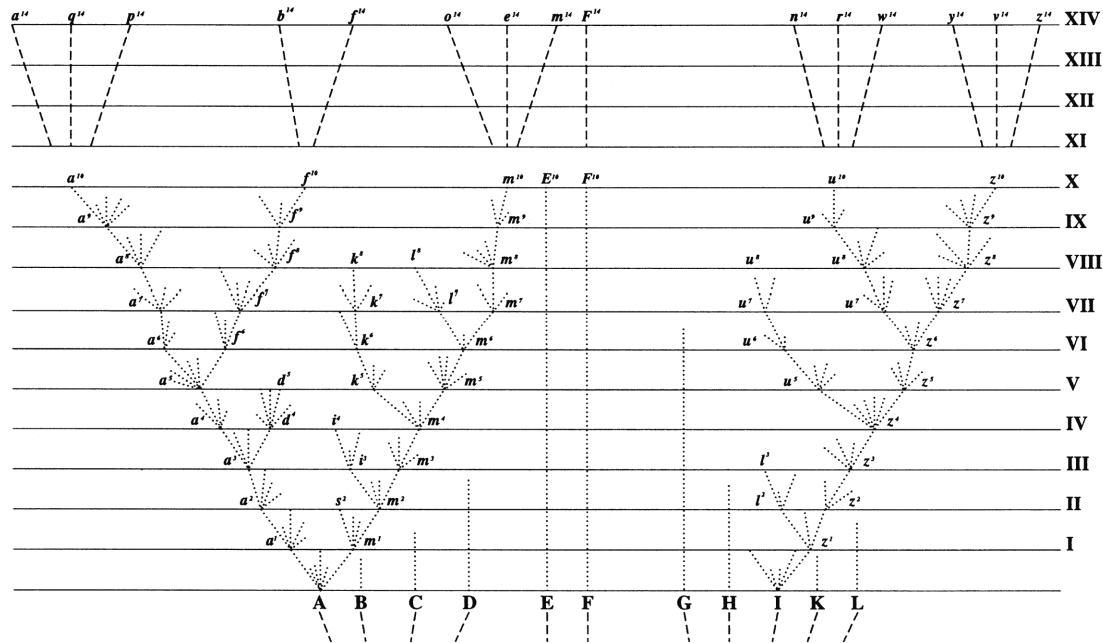
one more E be that one  
 lower than there are more  
 can more living on more  
 can more E be that one  
 as the same.

# CHARLES DARWIN QUOTE



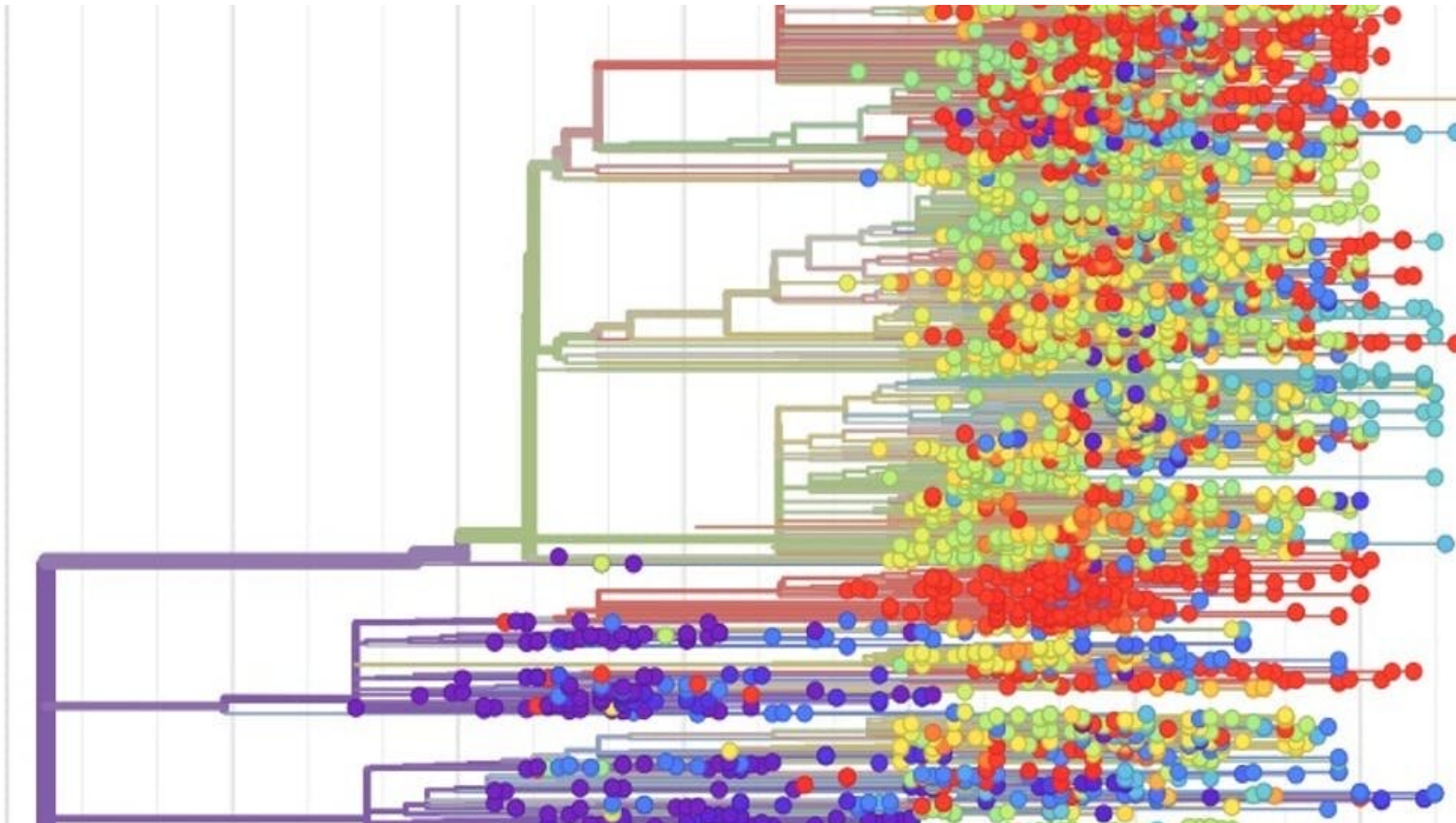
"It is not the strongest of the species that survives, nor the most intelligent that survives. It is the one that is most adaptable to change".

Charles Darwin





# SARS-CoV-2 (COVID) phylogenetic tree



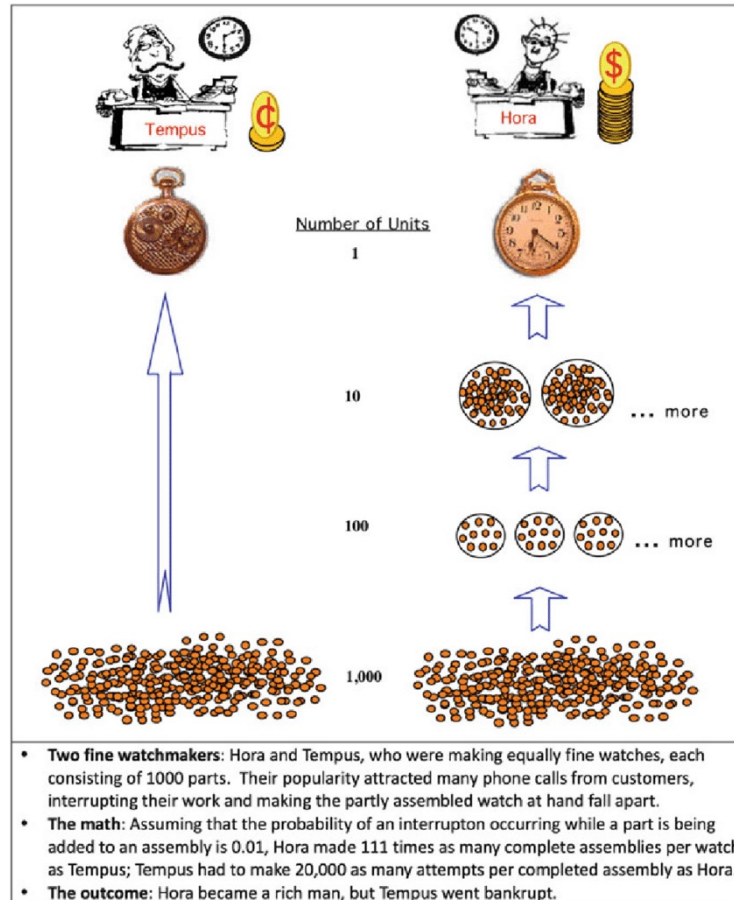
The SARS-CoV-2 phylogenetic tree – the family tree that shows the evolution of all the sequenced coronavirus samples worldwide.

# Herbert Simon's Parable on Evolution

## "The Parable of the Two Watchmakers"

A mathematical theory of "interruptions"

How can we quantify the Speed of Evolution?



HORA

TEMPUS

# Chapter 4: Hidden Markov Models



Image, courtesy of Vincent van Gogh Museum



# Chapter 4

---

## Machine Learning

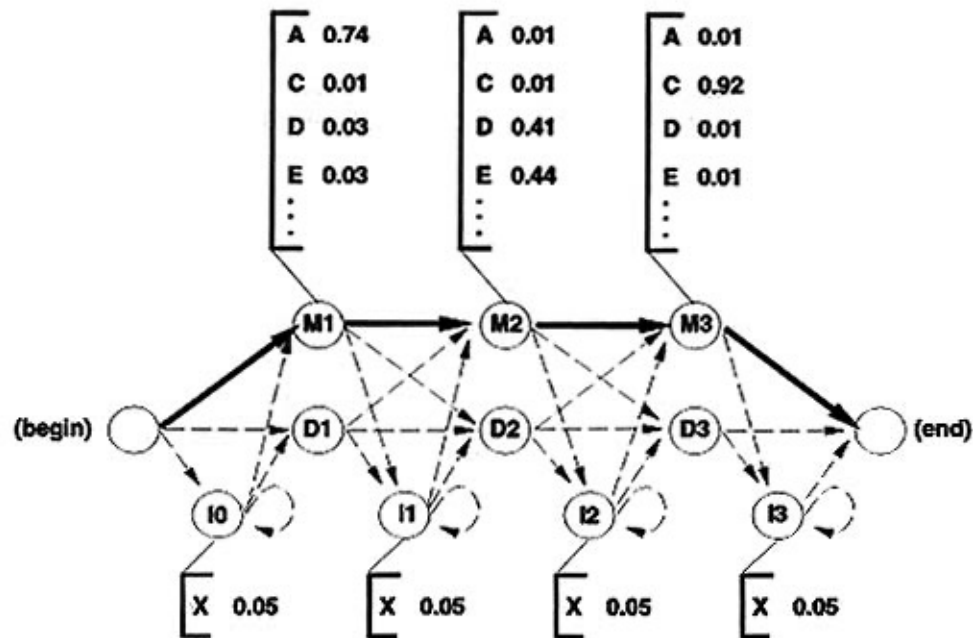
### Methods:

### Hidden Markov Models

### Algorithms

# Chapter 4: Hidden Markov Chains

## Algorithms



# Gene finding in a genome using HMMs algorithms

---



“For one rational line or true sentence there are thousands of nonsense cacophonies, mountains of verbal trash and incoherencies.”

Jorge Luis Borges

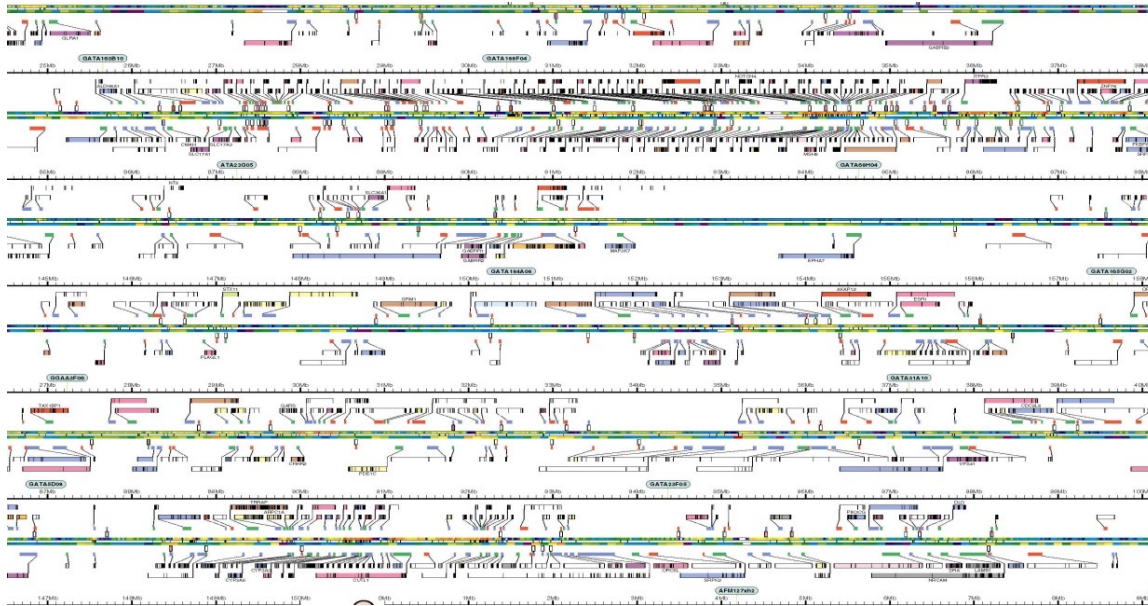
# Chapter 5: Genome Assembly Algorithms: An Introduction

---



# Chapter 5: Genome Assembly Algorithms

---

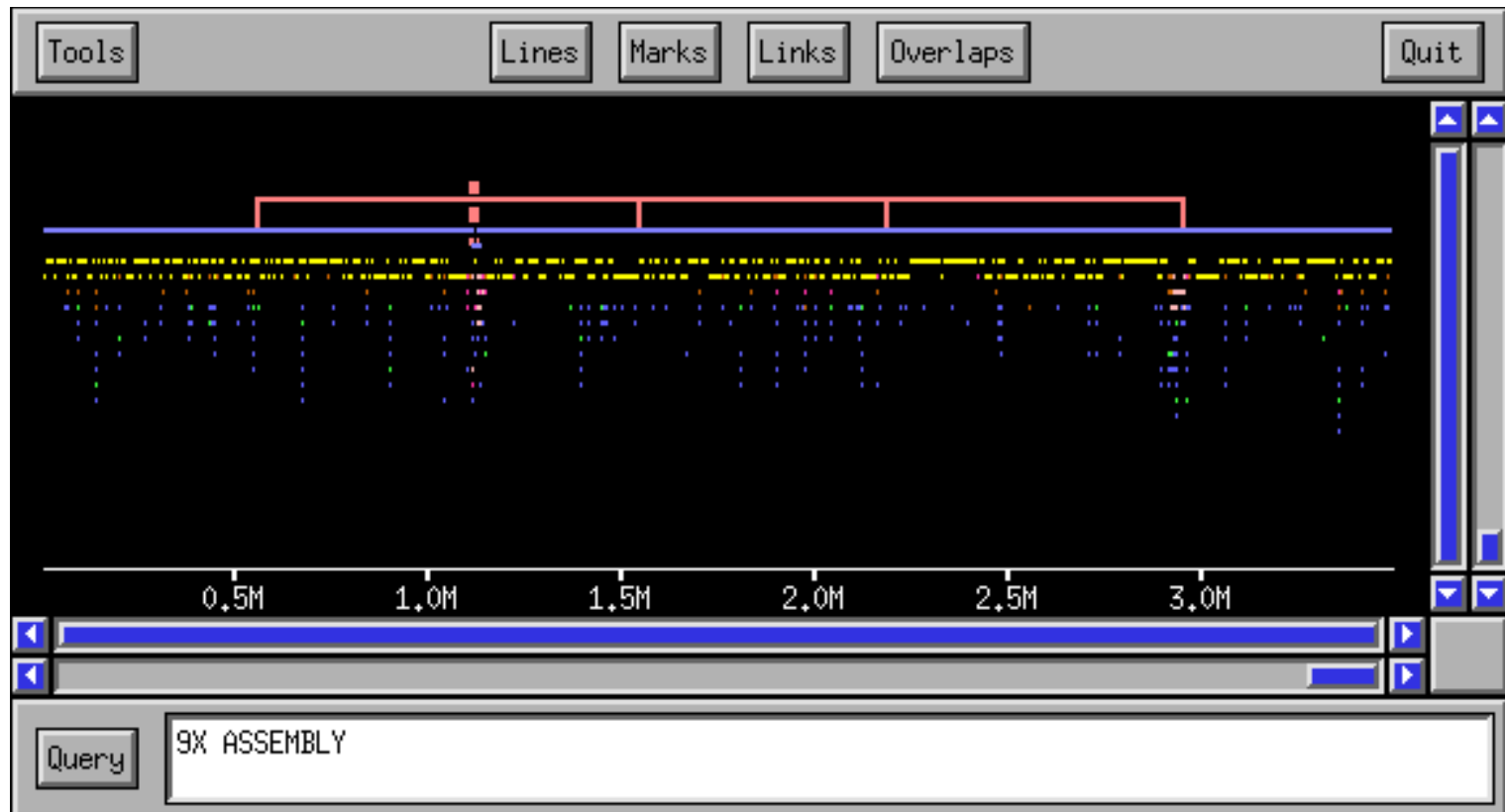






# Genome Assembly Algorithm

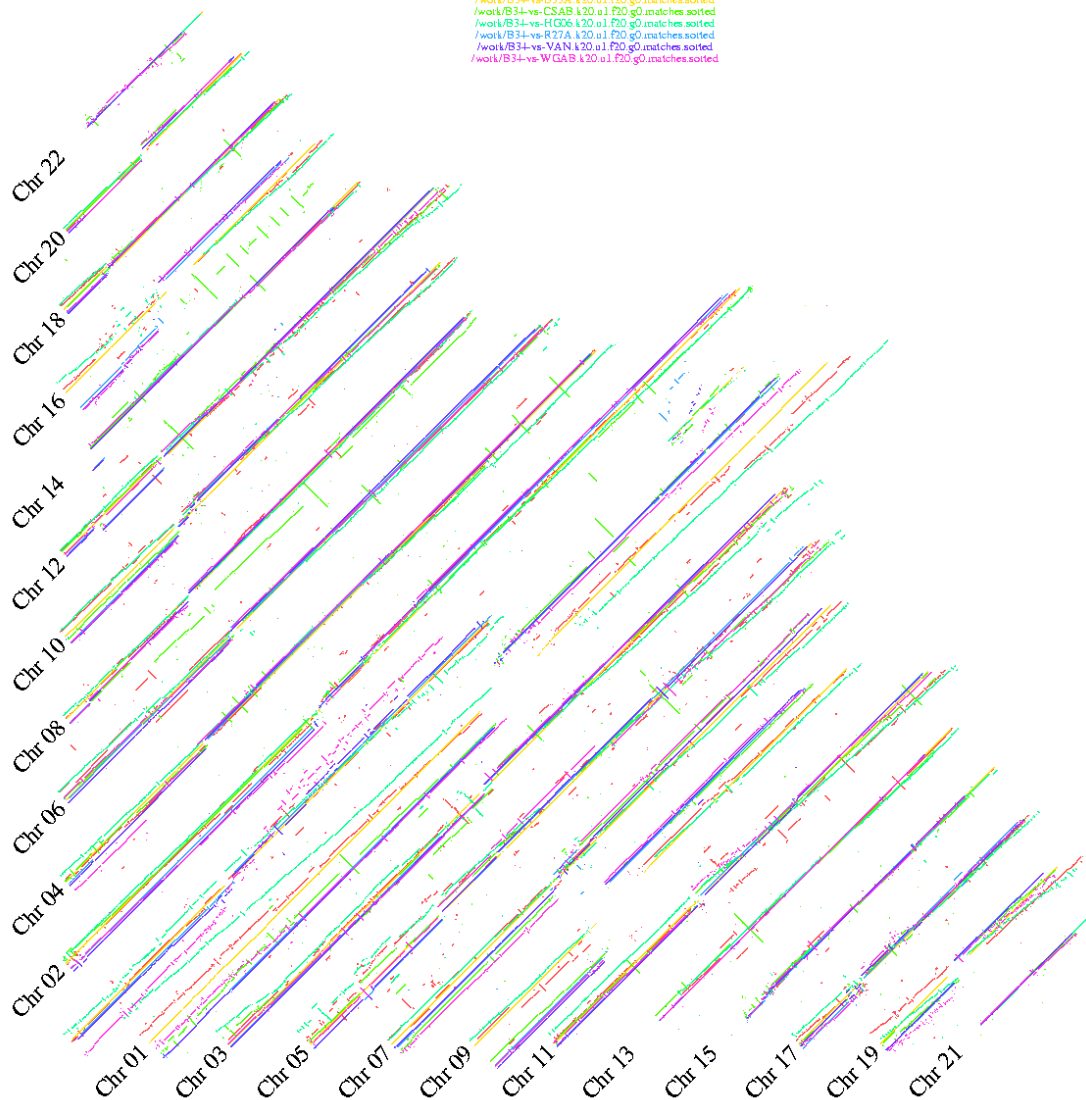
## Celera Assembler



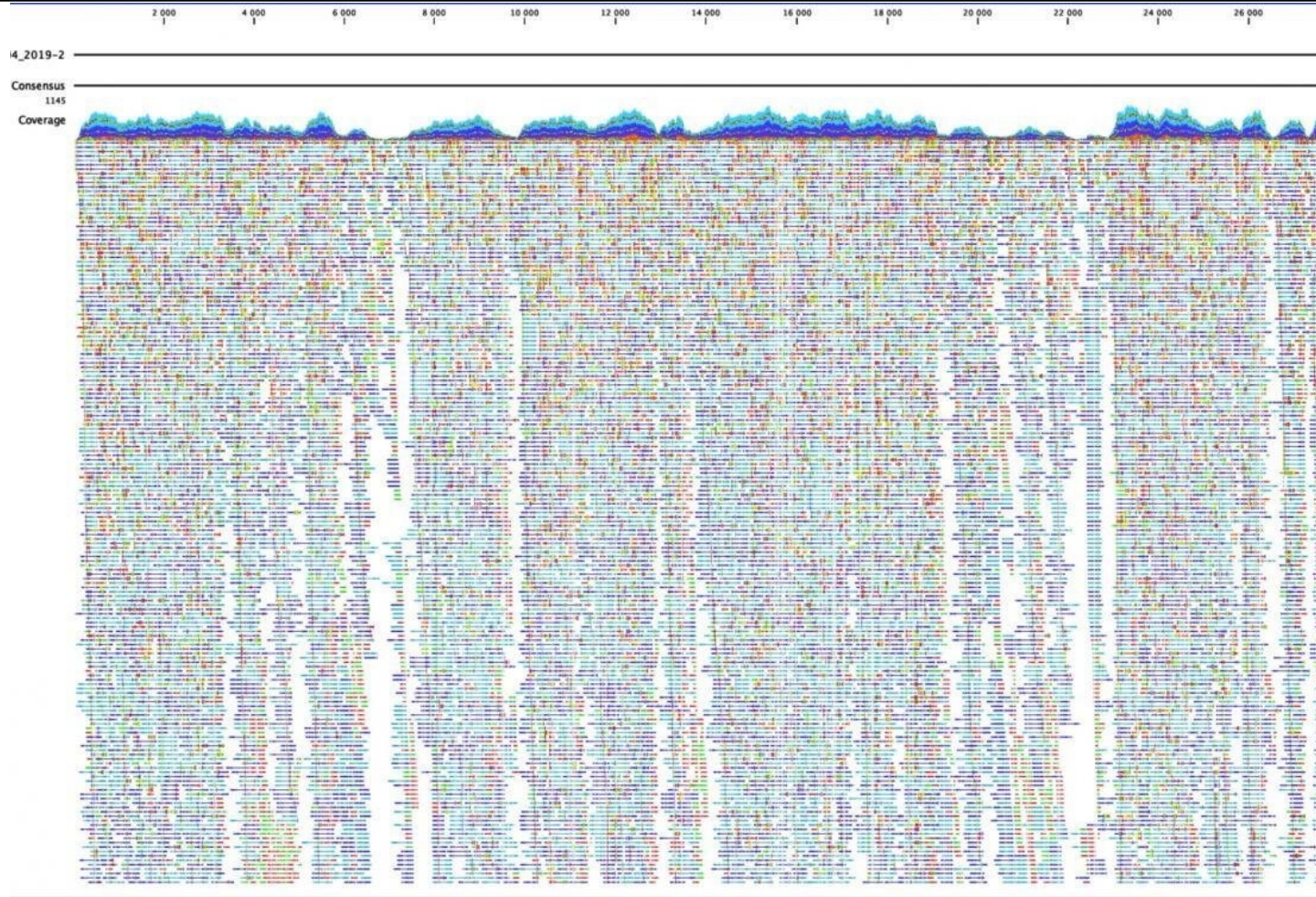


# The Father of All Dot Plots

```
/work/B34-vs-B2&.k20.u1.f20.g0.matches.sorted  
/work/B34-vs-B33A.k20.u1.f20.g0.matches.sorted  
/work/B34-vs-CSAB.k20.u1.f20.g0.matches.sorted  
/work/B34-vs-HCG6.k20.u1.f20.g0.matches.sorted  
/work/B34-vs-R27A.k20.u1.f20.g0.matches.sorted  
/work/B34-vs-VAN.k20.u1.f20.g0.matches.sorted  
/work/B34-vs-WCAB.k20.u1.f20.g0.matches.sorted
```



# The Human Genome



Whole genome sequence of the 2019-nCoV **coronavirus**, in one of the first French cases, made at the Institut Pasteur (Paris), using a unique Platform (P2M), open to all French National Reference Centers. Credit: Institut Pasteur/CNR of respiratory infection viruses.

# Chapter 6: Genomic Privacy





### **HOMER's attack:**

Genomic privacy studies on Genome-Wide Association Studies (GWAS) were first introduced as the well-known Homer's attack (2008) that showed that publicly released GWAS statistics can be used to estimate a GWAS participant's disease status from knowing her/his genotypes at certain risk factors.



**Do not**

**make the same mistake as Homer**

Take **CSCI2820: Medical Bioinformatics** and study the genetics of complex disease.

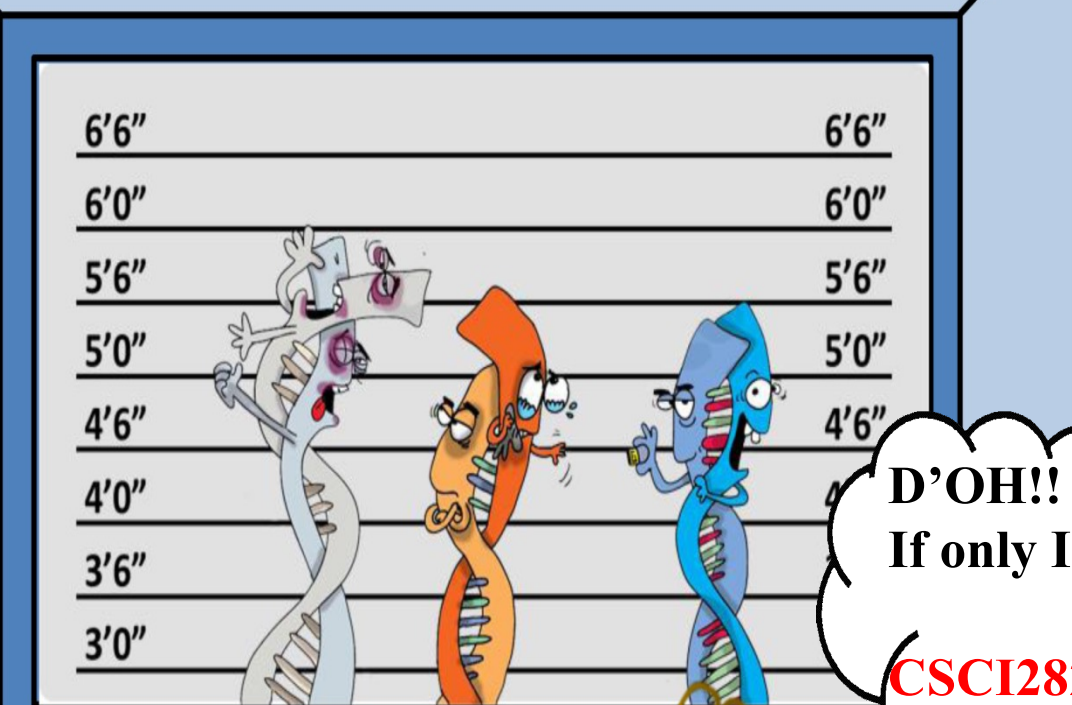
**Topics include**

- Hypothesis testing
- Haplotype phasing
- The missing heritability problem
- Genome-wide Association Studies
- Tag SNP selection
- Coalescent Theory

**CSCI2820**

**Tuesday/Thursday 2:30-3:50**

**Cartoon by  
prof. Alper Uzun (Brown Medical School)**



**D'OH!!  
If only I took  
CSCI2820...**

**Homer,  
implicate the  
gene that  
made you  
bald!**

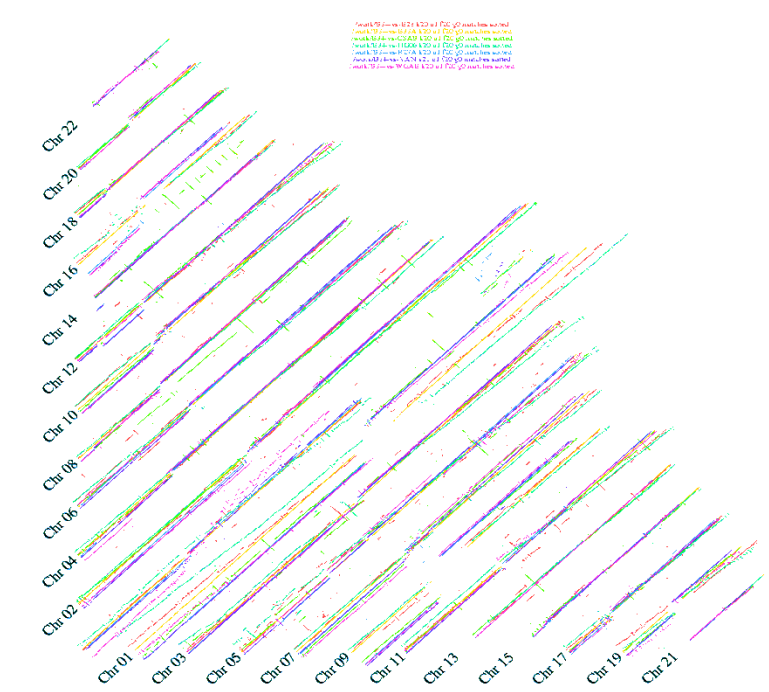




# CSCI1820 Algorithmic Foundations of Computational Biology

<http://www.cs.brown.edu/courses/csci1820/>

Prof. Sorin Istrail



"Whole-genome shotgun assembly and comparison of human genome assemblies" Proc. Nat. Acad. Sci. USA, 2004

"The Sequence of the Human Genome" Science, 2001

# Science

16 February 2001

Vol. 291 No. 5507  
Pages 1145-1434 \$9

## THE HUMAN GENOME

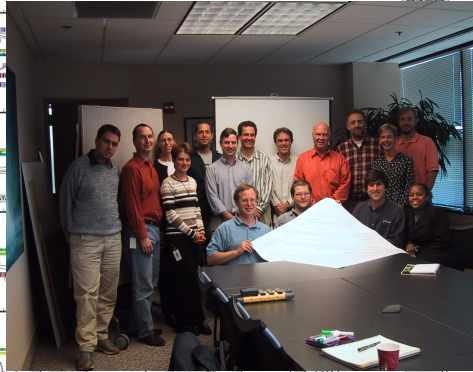


 AMERICAN ASSOCIATION FOR THE ADVANCEMENT OF SCIENCE

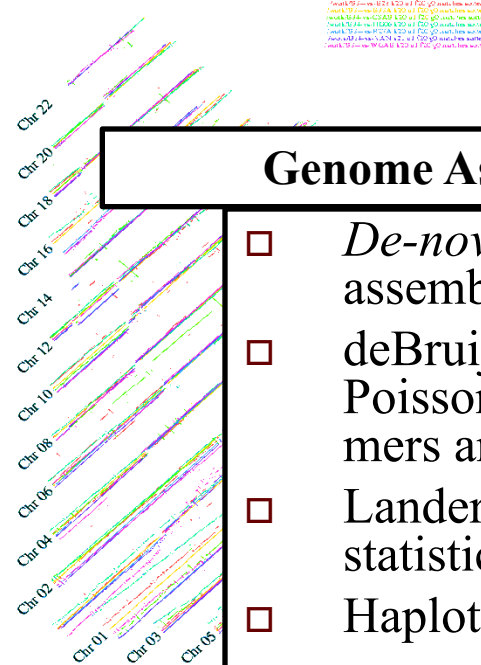
# CSCI1820 Algorithmic Foundations of Computational Biology

<http://www.cs.brown.edu/courses/csci1820/>

Prof. Sorin Istrail



"The Sequence of the Human Genome" Science, 2001



## Genome Assembly

- *De-novo* genome assembly algorithms
- deBruijn graphs and Poisson theory of k-mers and NGS
- Lander-Waterman statistical theory
- Haplotype assembly

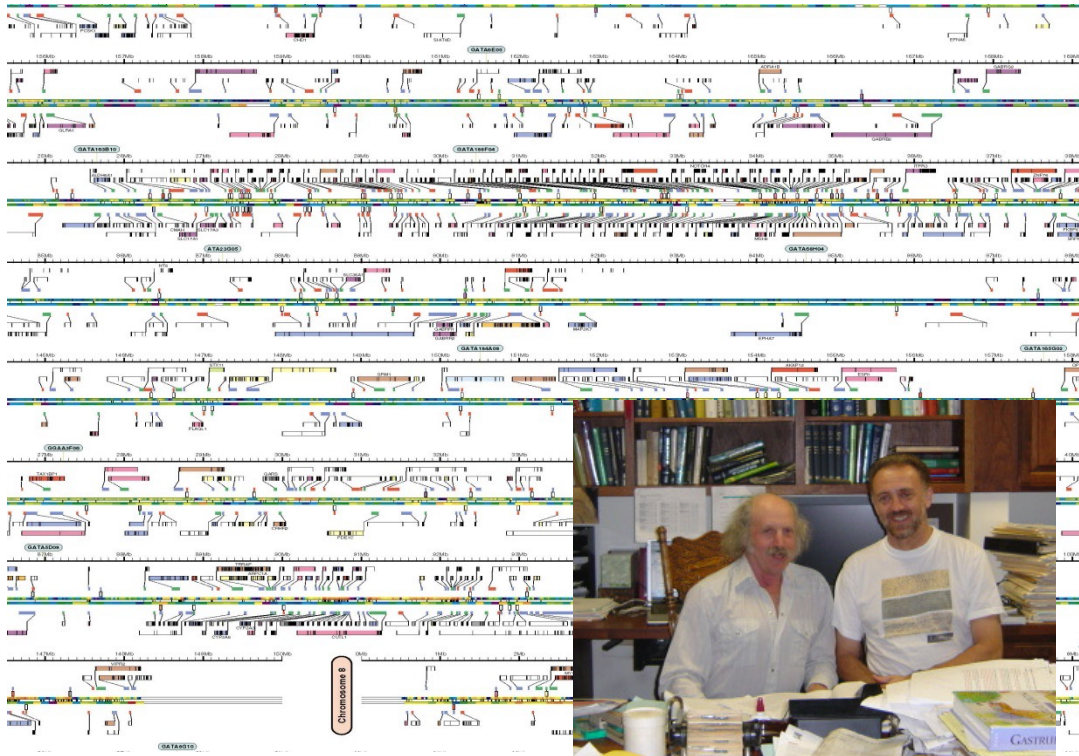
"Whole-genome shotgun assembly and comparison of human genome assemblies" Proc. Nat. Acad. Sci. USA, 2004

# CSCI1820 Algorithmic Foundations of Computational Biology

<http://www.cs.brown.edu/courses/csci1820/>

Prof. Sorin Istrail

The Regulatory Genome



Property	Value
CDS Links	14
Gene Accession	hCG181
Transcript Acc...	hCT195
Protein Acces...	hCP175
JAM Link	http://jss
Blast: Nasa Pr...	http://jss
Id	CELERA
Order Number	0
Comments	
Curator Flags	0
Feature Type	Transcript
Algorithm: Data	Curator
Parent Feature Id	hMORFS
Axis Name	GA_19F
Axis Id	CELERA
Axis Begin	294125
Axis End	324897

- The social network of Transcription Factors
- DNA combinatorics and statistics of Transcription Factor regulatory architecture
- Poisson clumps heuristics
- Chen-Stein Statistics

The cis-Regulatory CYRENE Genome Browser

Eric Davidson and Sorin working on "Logic functions of the genetic cis-regulatory code"

# CSCI1820 Algorithmic Foundations of Computational Biology

<http://www.cs.brown.edu/courses/csci1820/>

Prof. Sorin Istrail

- Suffix trees in linear time
- Burrows-Wheeler transform
- Karlin-Altschuler Statistics
- BLAST algorithm: Random walks, Information theory and P-values



CYRENE C

File Edit Search

onylocentrotus purpurus

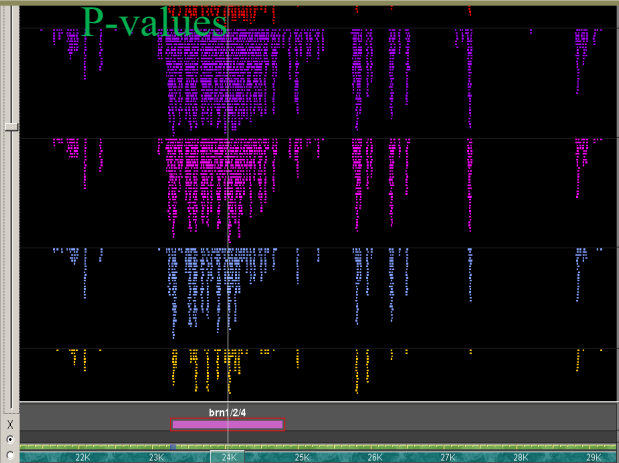
Unknown Chromosome

Genomic Axis (Left)

---

Promoted: RYAN:20071212

Property	Value
Id	20071212.10123
Order Number	1



The science and art of mapping DNA fragments, genes, and genomes to genomes

"The Sequence of the Human Genome" Science, 2001

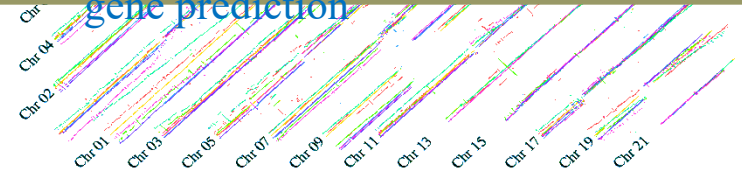
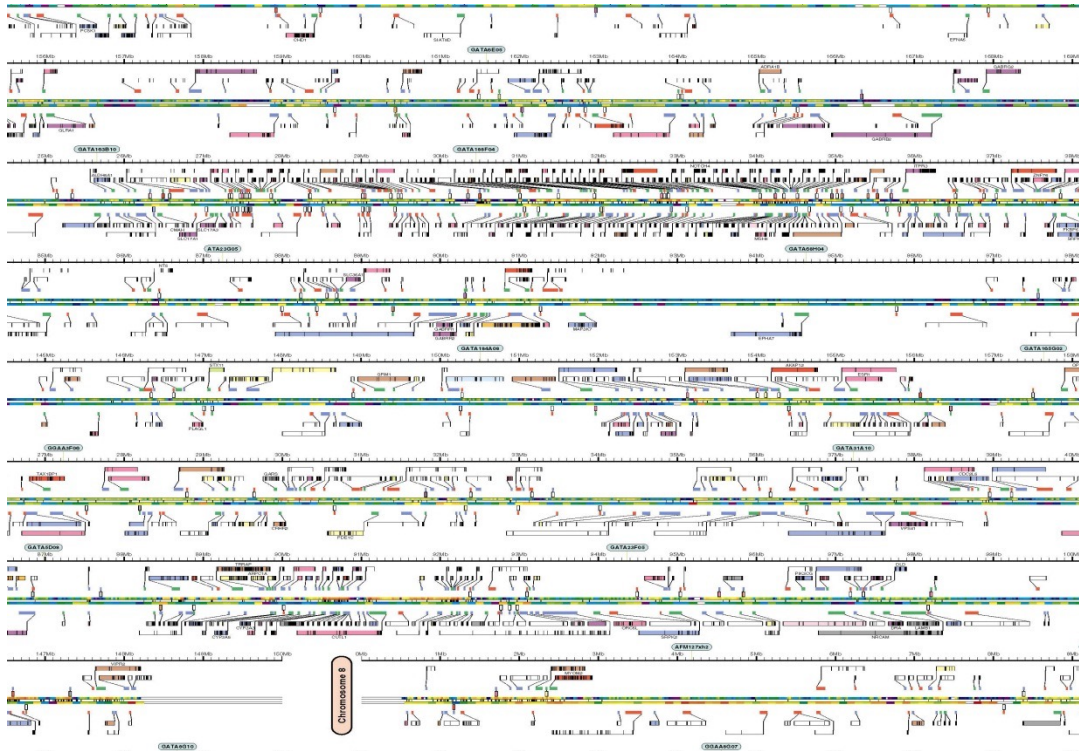
# CSCI1820 Algorithmic Foundations of Computational Biology

<http://www.cs.brown.edu/courses/csci1820/>

Prof. Sorin Istrail

Topics include

- Genome sequencing and assembly: algorithms and statistical theory
- BLAST algorithms and statistical theory of alignment and searching
- Mapping reads and genomes to genomes
- DNA combinatorics and statistical theory of regulatory regions of genes
- Hidden Markov model algorithms and gene prediction



"Whole-genome shotgun assembly and comparison of human genome assemblies" Proc. Nat. Acad. Sci. USA, 2004

"The Sequence of the Human Genome" Science, 2001



# CSCI2820 Advanced Algorithms in Computational Biology and Medical Bioinformatics

## Genome-wide Association Studies (GWAS)

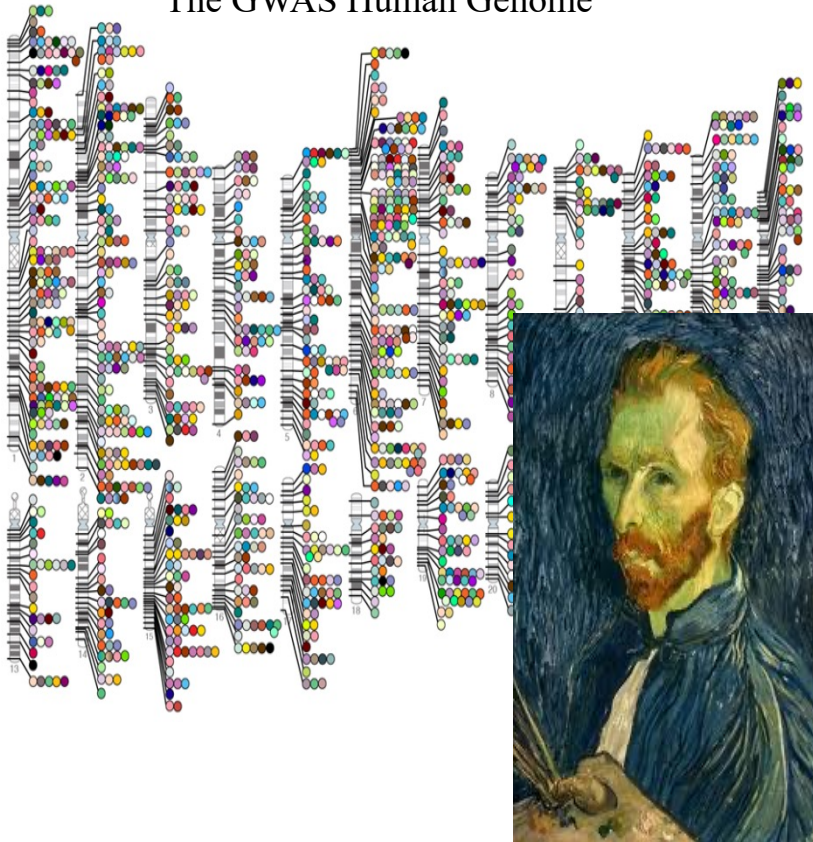
Prof. Sorin Istrail

Published Genome-Wide Associations through 2011

1,617 published GWA at  $p \leq 5 \times 10^{-8}$  for 249 traits



The GWAS Human Genome



### Genetic Heterogeneity

The Common Disease Common Variant (CDCV) hypothesis is dead.

Long live the Common Disease Many Rare Variants hypothesis!

The CDCV 's classical drawing metaphor as "Needles in the Haystack," with few needles with a common look in a large haystack, needs to be replaced now with a van Gogh-like drawing, with many needles each differently looking and private to areas in the large haystack.

Vincent



# CSCI2820 Advanced Algorithms in Computational Biology and Medical Bioinformatics

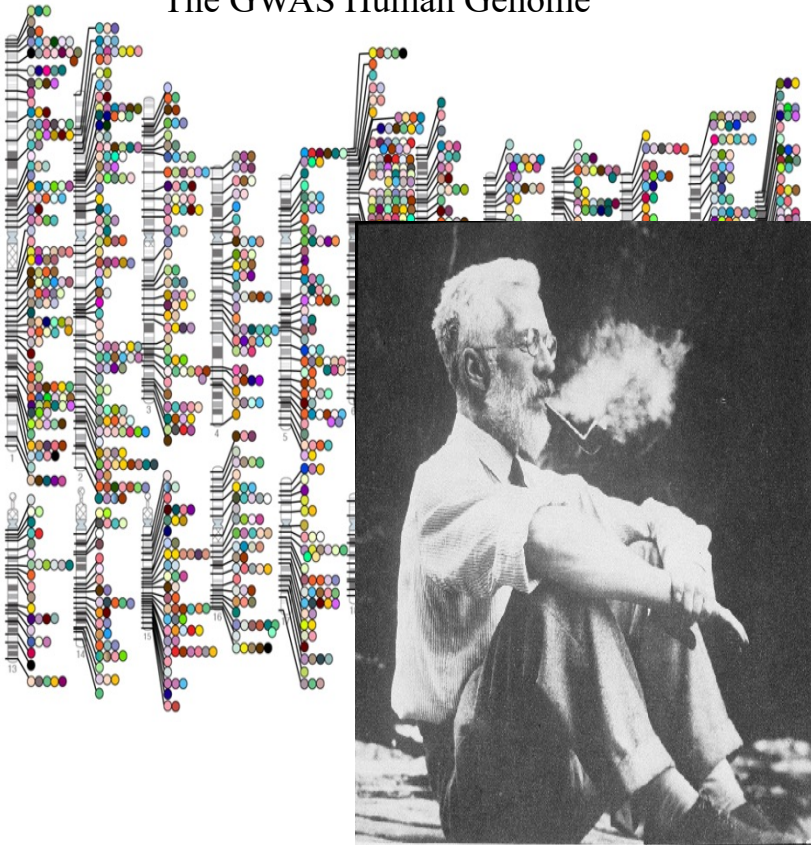
## Genome-wide Association Studies (GWAS)

Prof. Sorin Istrail

Published Genome-Wide Associations through 2011

1,617 published GWA at  $p \leq 5 \times 10^{-8}$  for 249 traits

The GWAS Human Genome



### The Missing Heritability Puzzle

Additivity of alleles? Just a convenient approximation, friendly to “heritability” measured as a correlation coefficient.

Ronald

# CSCI2820 Advanced Algorithms in Computational Biology and Medical Bioinformatics

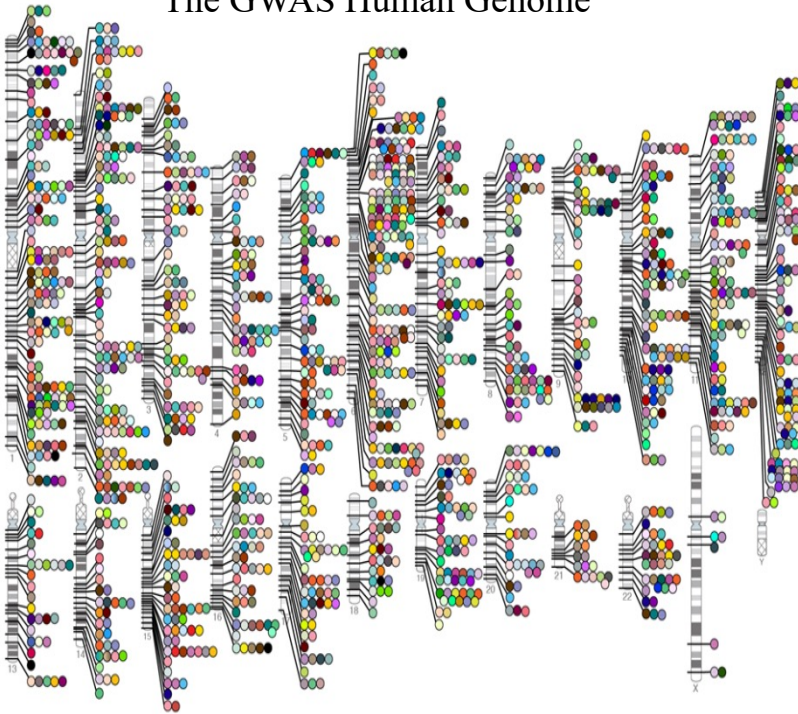
## Genome-wide Association Studies (GWAS)

Prof. Sorin Istrail

Published Genome-Wide Associations through 2011

1,617 published GWA at  $p \leq 5 \times 10^{-8}$  for 249 traits

The GWAS Human Genome



### Topics include

- haplotype phasing, linkage disequilibrium, tagging SNPs, identical by descent (IBD), pedigrees, trios
- coalescent theory, Polya urn game, Ewens sampling lemma, genome-wide graph theory algorithms
- the genetic heterogeneity problem, the missing heritability problem
- statistical models of disease, association tests and multiple hypothesis testing
- autism, multiple sclerosis, type 2 diabetes



# Bioinformatics is detective work

---

- **The Dancing Men code**, Sherlock Holmes
- **The Prison code**, a real life code used in CA

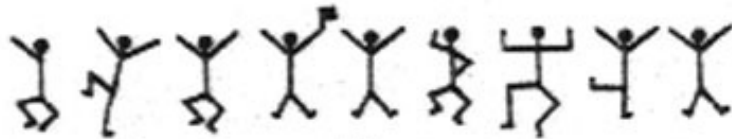
# The Adventures of the Dancing Men



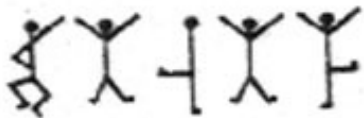
by Sir Arthur Conan Doyle  
"Sherlock Holmes"



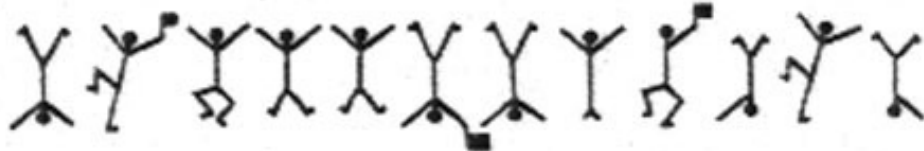
criminal's message (1)



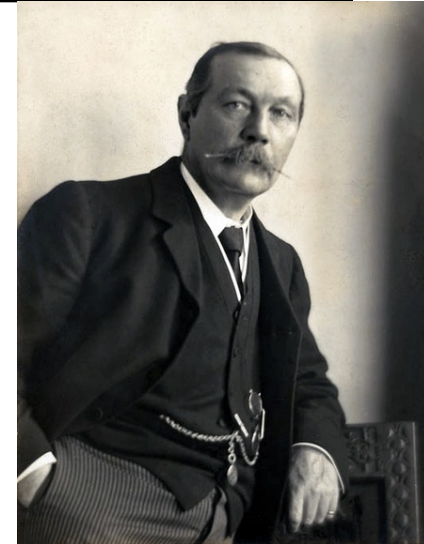
criminal's message (2)



Elsie's reply



criminal's message (3)



# The Adventures of the Dancing Men



by Sir Arthur Conan Doyle  
“Sherlock Holmes”

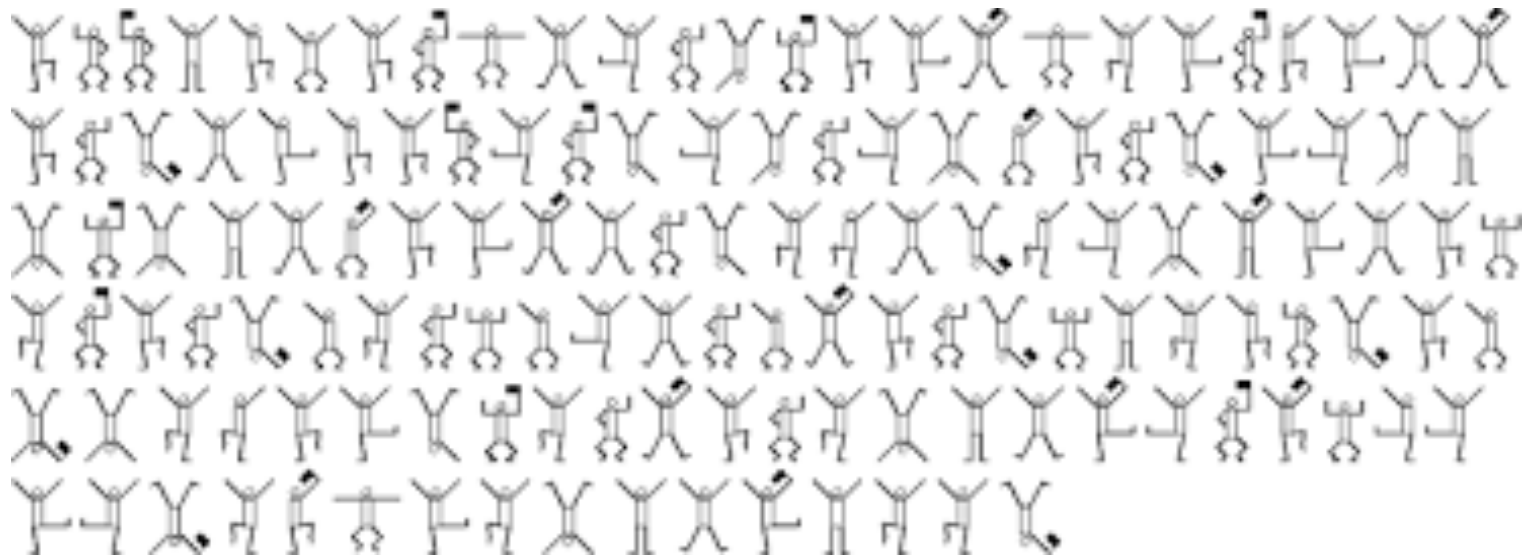
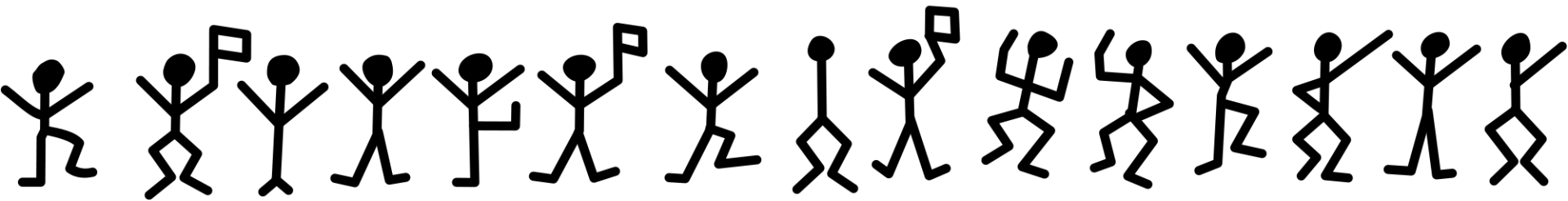
## Dancing Man Code (Sherlock Holmes)



A	B	C	D	E	F	G	H	I	J	K	L	M
N	O	P	Q	R	S	T	U	V	W	X	Y	Z
1	2	3	4	5	6	7	8	9	0			

# The Dancing Men code

by Arthur Conan Doyle: "The Adventures of the Dancing Men"



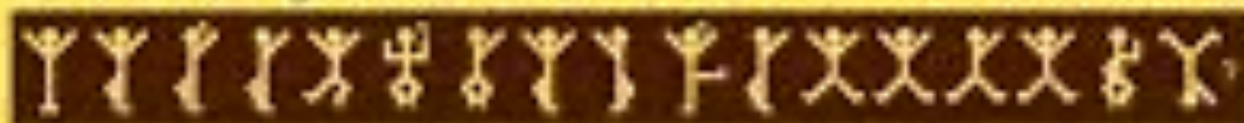
Hello, how are you?



Good morning!



How was your weekend?



























What the heck!



I've seen that movie!



 E	 T	 A	 O	 I/J	 R
 S	 U	 N	 F	 L	 P
 H	 B	 D	 G	 W	 V
 M	 K	 Y	 C	 X	 Q/Z



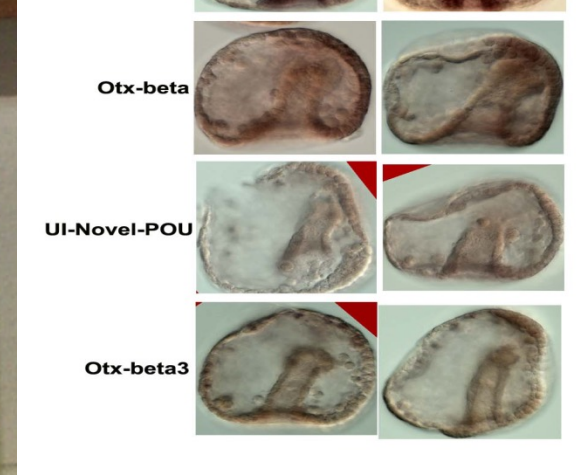


Handwritten text in a stylized, cursive script, possibly representing a prison code or a specific dialect. The text is arranged in several lines across a grid of horizontal lines.



# The Prison code

**Solution:** An Algorithm based on Markov Chain Monte Carlo



Caltech, Davidson Lab  
October 2004



Eric Davidson  
– in memoriam