

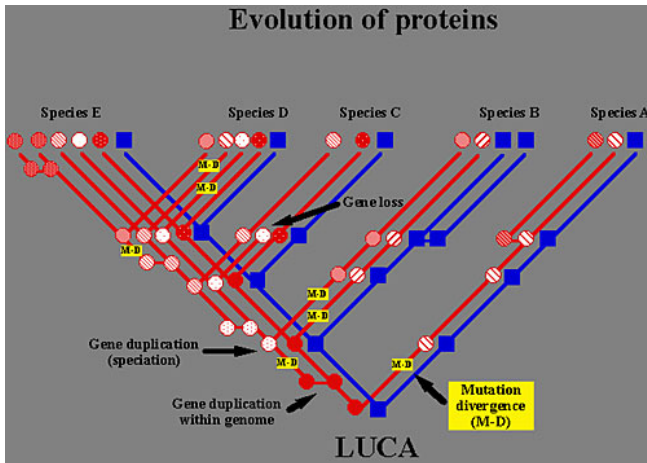
Computational Biology

What for ?

Protein family

Gene sequence

Structure



```

420      *      440      *      460      *
Hs-RPL3-I3 : AATGGGTGTAACCTAA-ATTAACCTGTGGACCTCTGCTCAGCTCCGCTCGGCTCTGCC : 45
Mm-RPL3-I3 : TATATTTGGCTAT-ATTAACCT-GTGGACCTCTGCTCAGCTCCGCTCGGCTCTGCC : 42
Bt-RPL3-I3 : GAATGTGTGTTTCATTATTAACCTGTGGACCTCTGCTCAGCTCCGCTCGGCTCTGCC : 39
At tTGTg cTtAt ATTAACCTGTGGACCTCTGCTCAGCTCCGCTCGGCTCTGCC

480      *      500      *      520      *
Hs-RPL3-I3 : CGATGAGCTCCATCCAGGCTCCGCTGCGCGTGGAAAAGGCTCCTTAGAAGCCGGCAAT : 51
Mm-RPL3-I3 : CGATGAGCTTCATCCAGGCTCCGCTGCGCGTGGAAAAGGCTCCTTAGAAGCCGGCAAT : 47
Bt-RPL3-I3 : CGATGAGCTCCATCCAGGCTCCGCTGCGCGTGGAAAAGGCTCCTTAGAAGCCGGCAAT : 45
CGATGAGCTcCATCCAGGCTCCGCTTGCCTGGAAAAGGCTCCTTAGAAGCCGGCAAT

540      *      560      *      580      *
Hs-RPL3-I3 : GAGCTCCATCCCCACGGTGCAGTGTGCCTTCCGCTCACCCCTCGAGGGGTGATGA : 57
Mm-RPL3-I3 : GAGCCCCATCCCCAATGGTGCAGTGTGCCTTCCCTCACCTGTTGCAAGGTTGATGA : 53
Bt-RPL3-I3 : GAGCCCCCTCCCCACACGGTGCCTGTGTGCCTCCCCCTCACCCGTTGGAGGGTGTGATGA : 50
GAGCCCaTCCCCAcacGGTGCcAgTgTGCcTtCc CTCACcGtTtGgAgGgGtGATGA

600      *      620      *
Hs-RPL3-I3 : AGGCCTGCACC-TGGTCCCCTCCCCAACTCTGCTCTGCTCCTGAAG : 619
Mm-RPL3-I3 : AGGCCTGCACC-GGGCCCCTCCCCAACTCTGCTCTGCTCCTGAAG : 583
Bt-RPL3-I3 : AGGCTGGCACCTGGGCCCCTCCCCAACTCTGCTCTGCTCCTGAAG : 555
AGGCcTGCACC gGgcccctCCCCAACTcTgCTCTgCTcCTGAAG
    
```

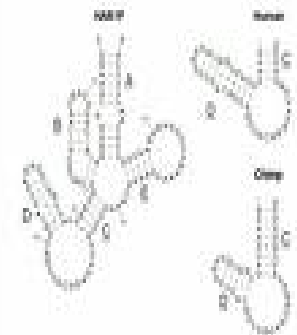
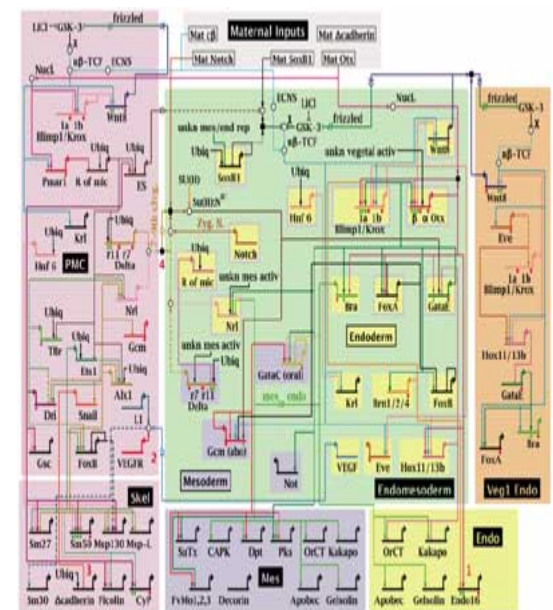
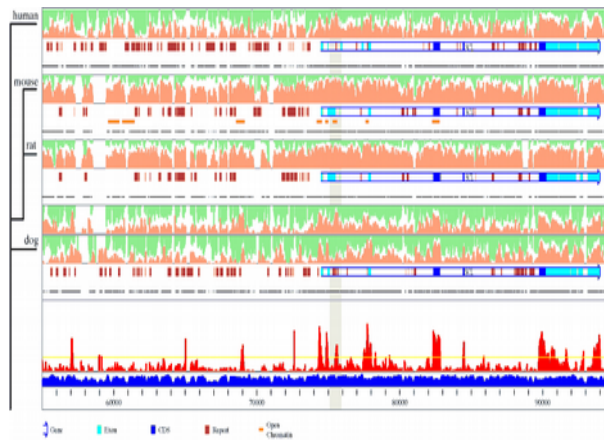
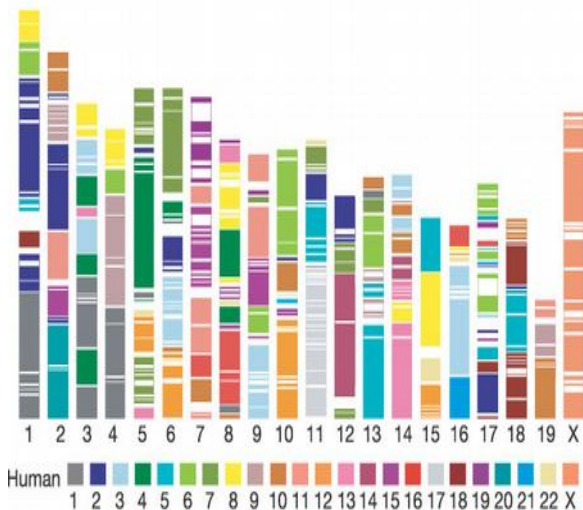


Figure 2. Alignment of human, bovine and mouse RPL3 gene (intron 3) using the Genedoc multiple alignment program. Note the identical positions of nucleotides.

Gene order

Non-coding regions

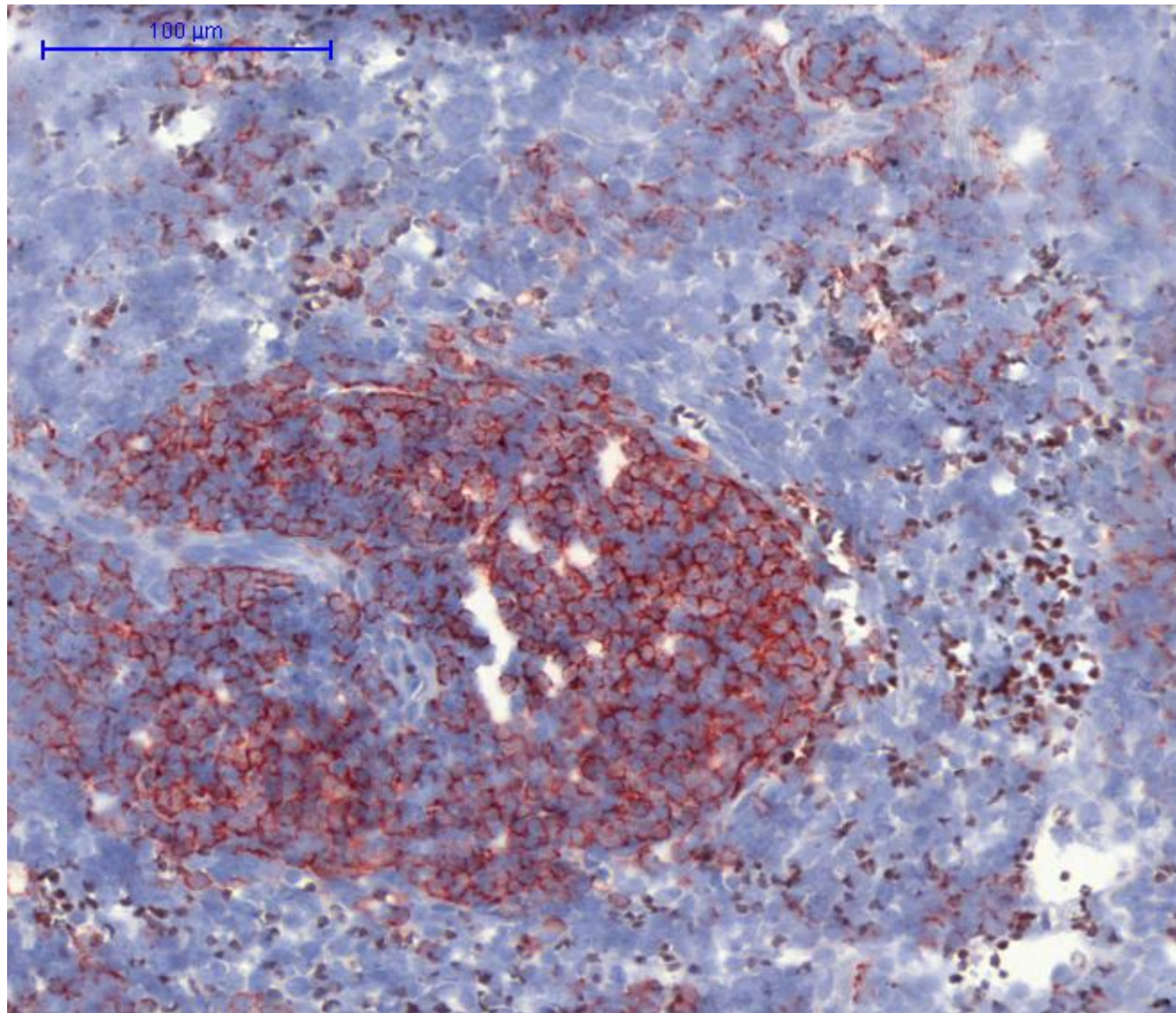
Gene networks



Computational Biology

What for ?

And their phenotypic counterpart



Translational Bioinformatics

PLOS Computational Biology

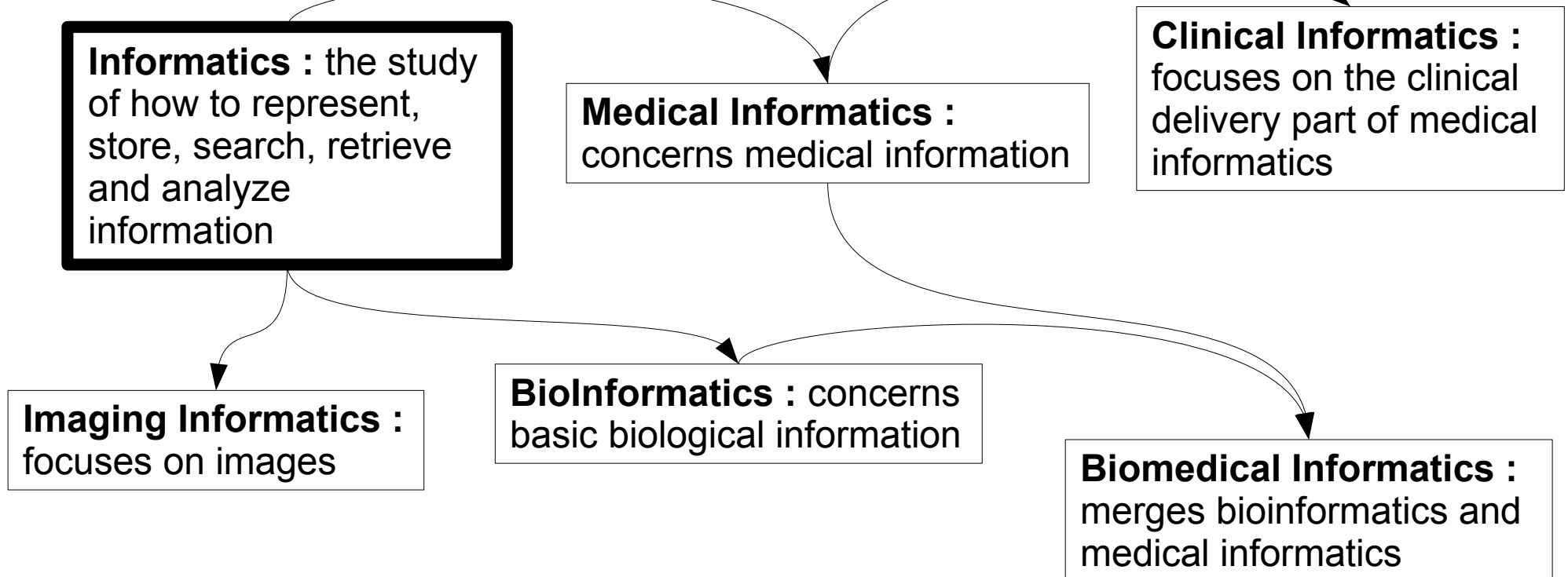
A Peer-Reviewed, Open Access Journal

Translational Bioinformatics : a collection of Education Articles, 2012

<http://www.ploscollections.org/article/browseIssue.action?issue=info:doi/10.1371/issue.pcol.v03.i11>

Impact of Computational Biology : translational sciences

Integrate huge amount of **heterogeneous** molecular and clinical data for a better understanding of molecular basis of diseases and subsequently **changing** clinical practices of course for the benefice of the patient



Translational :

- How to improve diagnostic, pronostic and patients' care ?
 - Small devices
 - Molecular dignostic
 - Nano-particules based treatment
 - Vaccine
 - Etc.
- Mastering the huge amount of new knowledge in molecular biology, genetics and genomic.

Double helixoidal structure of ADN → pratical improvement of human health from a technological point of view ?

For sure, we are able to quickly compute/measure :

- DNA sequences (whole genome scale)
- RNA sequences and expression
- protein sequences, structure, expression and modification
- structure, presence and quantity of small molecular metabolites
- generate a lot of data including images

2 playground chapters for this sessions :

- Quantitative Imagery
- Machine Learning / Data Mining

2 important chapters in an ideal world :

- Graph and Network representations
- Knowledge representations : data, database, ontologies

Then technologies/ environments for software use/development :

- Java : ImageJ, Weka
- Python : Biopython, Numpy, Scipy, Matplotlib, Pyvis, Enthought Python Distribution and Canopy, Anaconda Pandas
- scripts : Perl, Gawk
- Inkscape, ImageMagik / Sphinx / XML, SBML, BioPax, GPML, JSON, SQL, noSQL, Hadoop
- Clustal → T-coffee, PathwayAPI, BioGRID, PatternHunter

Learning code In a biological perspective for being able to get involved in...

- **Gene feature recognition :**

- **TIS (Translation Initiation Site)**
- **TSS (Transcriptional Start Sites)**
- **Feature Generation → Feature Selection → Feature integration**
- **Gene finding**

- **Gene expression analysis :**

- **Affymetrix Gene Chip Data**
- **Gene expression Profile Classification**
- **Gene expression Profile Clustering**
- **Gene Regulatory Circuits Reconstruction : Differentially Expressed Genes, Gene Interaction Prediction**

- **Sequence Alignment / Comparaison / Homology :**

- **Multiple Sequence Alignment (Dynamic Programming)**
- **Function assignment to protein sequence (Guilt by Association)**
- **Discovery of Active Site or Domain of a function**
- **PPI / Proteomic Profile Analysis**
- **key mutation site identification**

- **Phylogenetic tree :**

- **Construction**
- **Comparison**

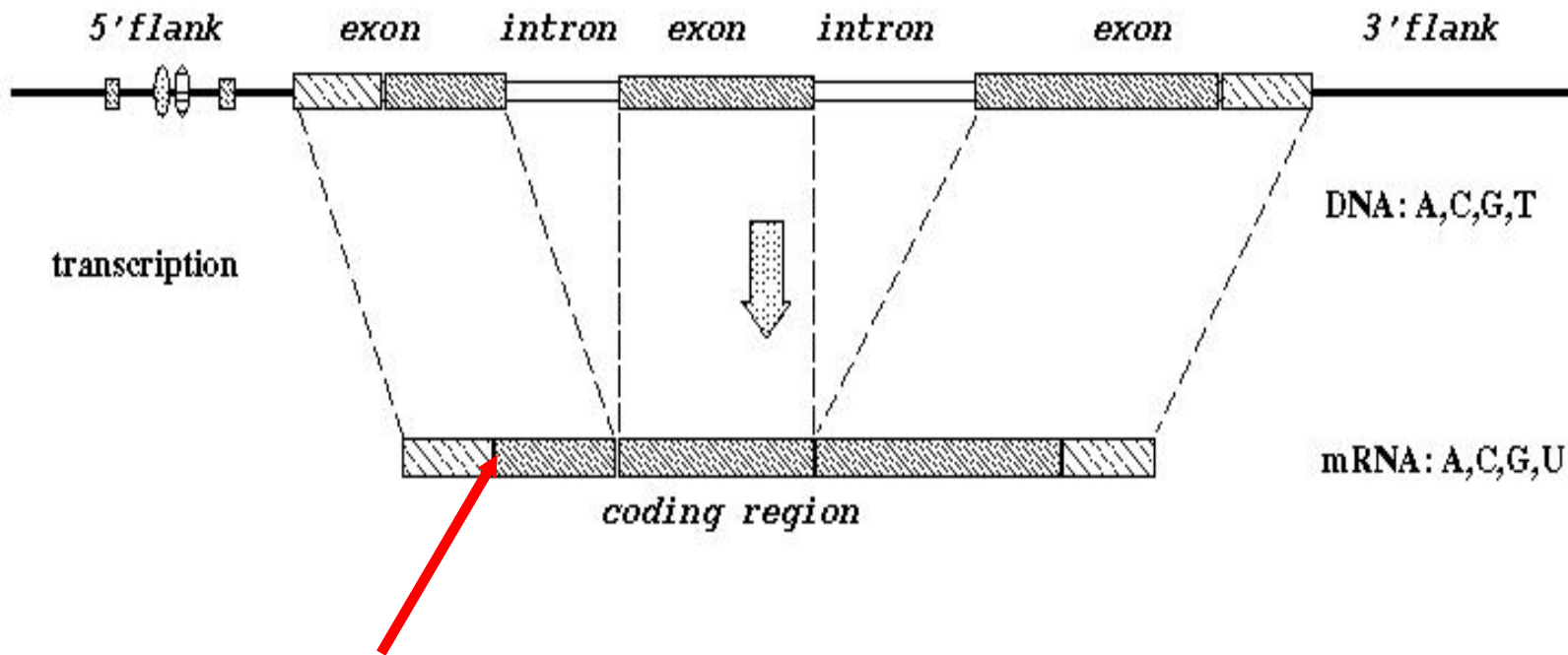
-**Biological Networks / Graph of Interactions :**

- **Natural pathways**
- **PPI Networks**
- **Protein Complex Prediction**

- **Image analysis :**

- **High-throughput screening**

TIS : Translation Initiation Site Recognition/Prediction



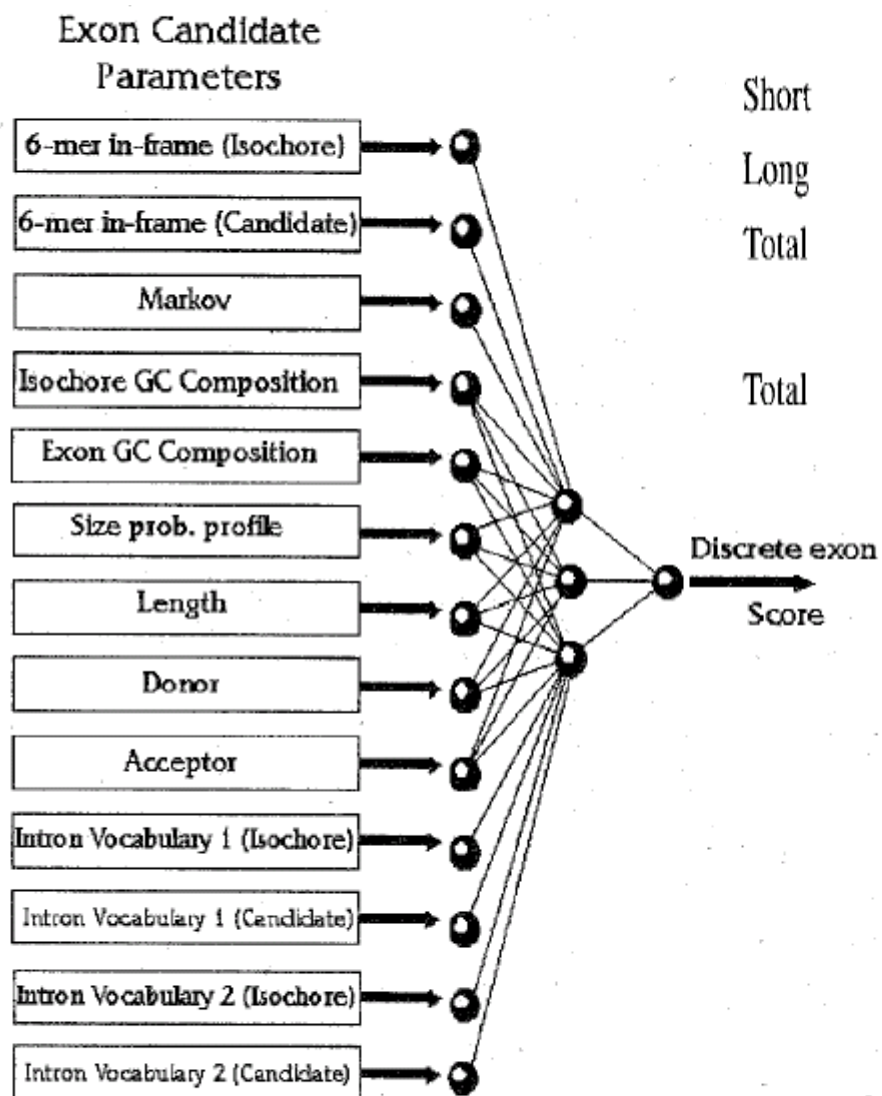
cDNA sample

```

299 HSU27655.1 CAT U27655 Homo sapiens
CGTGTGTGCAGCAGCCTGCAGCTGCCCAAGCCATGGCTGAACACTGACTCCCAGCTGTG      80
CCCAGGGCTTCAAAGACTTCTCAGCTTCGAGCATGGCTTTTGGCTGTCAGGGCAGCTGTA     160
GGAGGCAGATGAGAAGAGGGAGATGGCCTTGGAGGAAGGGAAGGGGCCTGGTGCCGAGGA     240
CCTCTCCTGGCCAGGAGCTTCTCCAGGACAAGACCTTCCACCCAACAAGGACTCCCCT
.....                                                                    80
.....iEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE     160
EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE     240
EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE
    
```

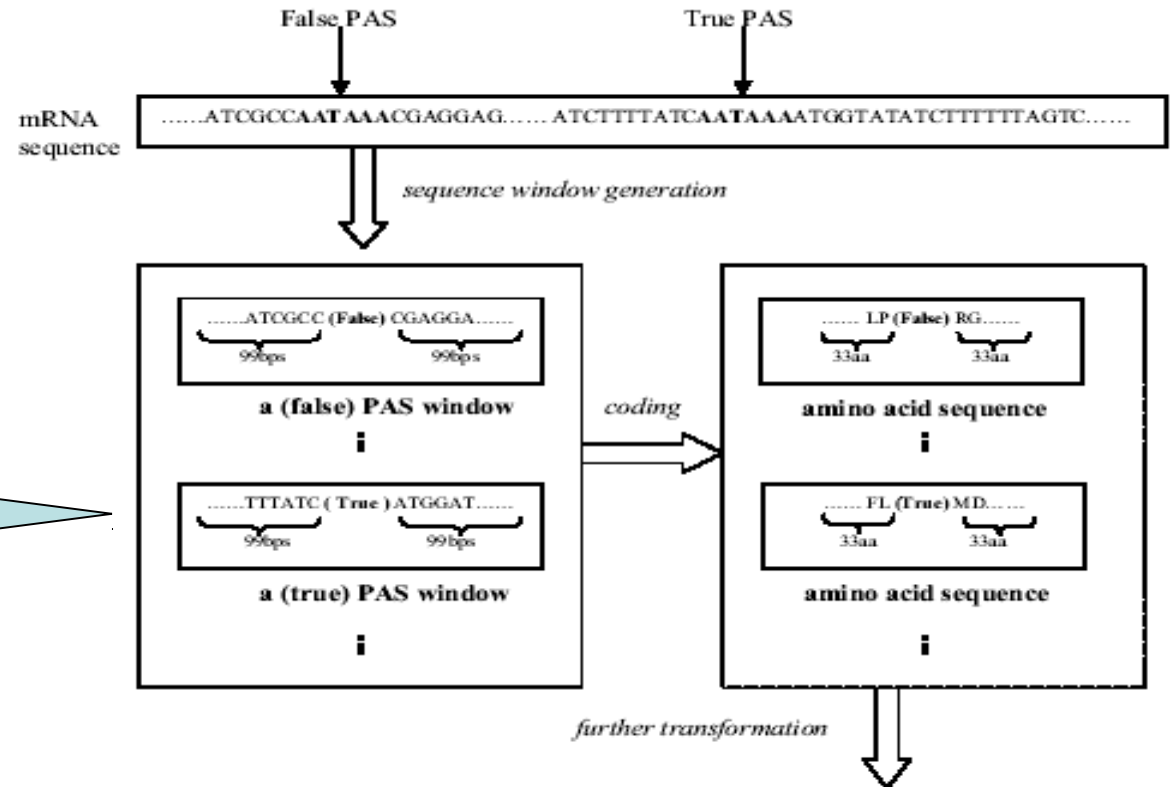
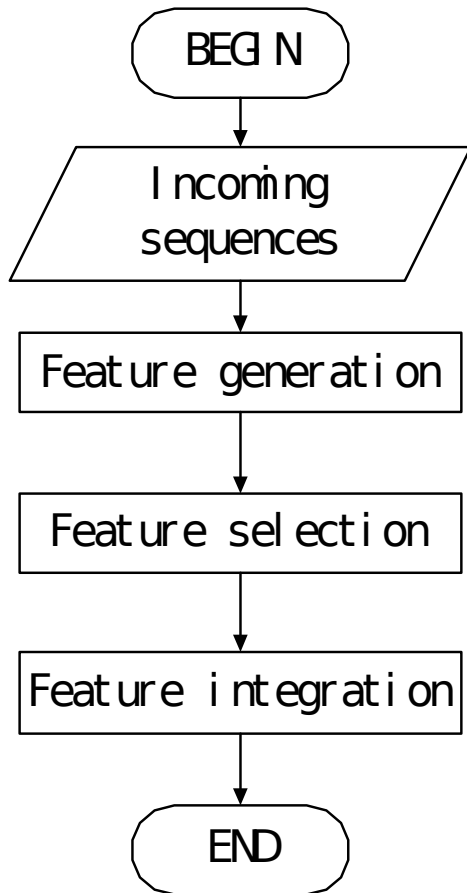
Why the second ATG is a TIS?

Gene prediction



DNA	Predictions			
	TP	%	FP	%
# Exons				
229	171	74.7	39	18.6
600	575	95.8	30	4.9
829	746	90.0	69	8.5
# Bases				
134814	122885	91.2	13048	9.6

PAS Prediction



New feature space (total of 925 features + class label)			
42 1-gram amino acid patterns	882 2-gram amino acid patterns	1 bio-knowledge pattern	class label
UP-A, UP-R, ..., UP-N, DOWN-A, DOWN-R, ..., DOWN-N (numeric type)	UP-AA, UP-AR, ..., UP-NN, DOWN-AA, DOWN-AR, ..., DOWN-NN (numeric type)	UP-T-number (numeric type)	True, False
Frequency as values			
1, 3, 5, 0, 4, ...	6, 2, 7, 0, 5, ...	10,	False
!	!	!	!
6, 5, 7, 9, 0, ...	2, 0, 3, 10, 0, ...	50,	True
⋮	⋮	⋮	⋮



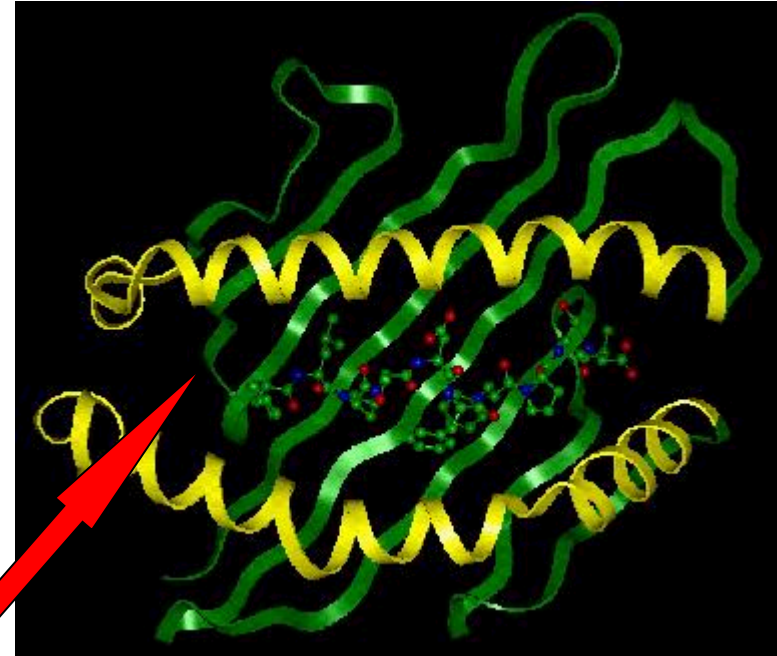
SVM in Weka

T-Cell Epitopes Prediction By Artificial Neural Network

- Honeyman et al., *Nature Biotechnology* 16:966-969, 1998

TRAP-559AA

MNHLGNVKYLVIVFLIFFDLFLVNGRDVQNNIVDEIKYSE
EVCNDQVDLYLLMDCSGSIRRHNVVNHAVPLAMKLIQQLN
LNDNAIHLY**VNVFSNNAK**EIIIRLHSDASKNKEKALIIIRS
LLSTNLPYGRTNLTDALLQVRKHLNDRINRENANQLVIL
TDGIPDSIQDSLKESRKLSDRGVKIAVFGIGQGGINVAFNR
FLVGCHPSDGKCNLYADSAWENV**KNVIGPFMKAVCVEVEK**
TASCGVWDEWSPCSVTGKGRSRKREILHEGCTSEIQEQ
CEEERCPPKWEPLDVPDEPEDDQPRP**RGDNS SVQK**PEENI
IDNNPQEPSPNPEEGKDENPNGFDLDENPENPPNPDIPF
KPNIPEDSEKEVPSDVPKNPEDDREENFDIPKKPENKDN
QNNLPNDKSDRN**IPYSPLPPK**VLDNERKQSDPQSDMNGN
RHVPNSDRETRPHGRNNENRSYNRKYNDTPKHFLREEHE
KPDNNKKKGESDNKYKIAGGIAGGLAL**LACAGLAYK**FVVP
GAATPYAGEPAPFDETLGEEDKDLDEPEQFRLPEENEWN



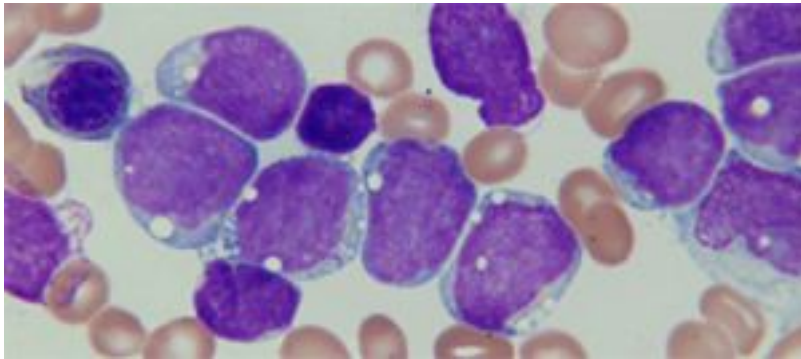
Histone Promoter Recognition Programs

“GENERAL PROMOTERS”			
First Generation			
Name	Scoring Technique used	Search by content/signal	Features used
NNPP	Time delay NN	Signal	TATA box, Inr
Promoter 2	NN	Signal	TATA box, Inr, CAAT box, GC Box
PromFind	Discriminative count	Content	Hexamer frequency
PromoterScan	Discriminative count	Signal	TATA box, TFBS
TSSG/TSSW	Linear discriminant analysis	Content + signal	TATA box, TSS, hexamer frequency upstream TSS, TFBS
Second Generation			
DGSF	NN	Content + signal	CpG island, TSS, DPF output
DPF	NN	Content + signal	Promoter, exon, intron, TSS
Eponine	SVM variant	Content + signal	TATA box, GC rich content, TSS
FirstEF	Quadratic discriminant analysis	Content + signal	First exon, CpG islands
Mcpromoter	NN & Interpolated markov models	Content + signal	TATA box, CAAT box, GC box, nucleosome position
PromoterInspector	Discriminative counts	Content	Oligonucleotides, Exon, Intron, 3'UTR, Promoter genomic context
CpG Promoter	Quadratic discriminant analysis	Content + signal	CpG island, TSS
CpGProD	Generalised linear model	Content	CpG island, AT/GC content
“SUB-CLASS OF PROMOTERS”			
Muscle family	Discriminative counts	Signal	TFBS, relative distance
Globin family	Logical operators AND, OR and NOT	Signal	TFBS, relative distance

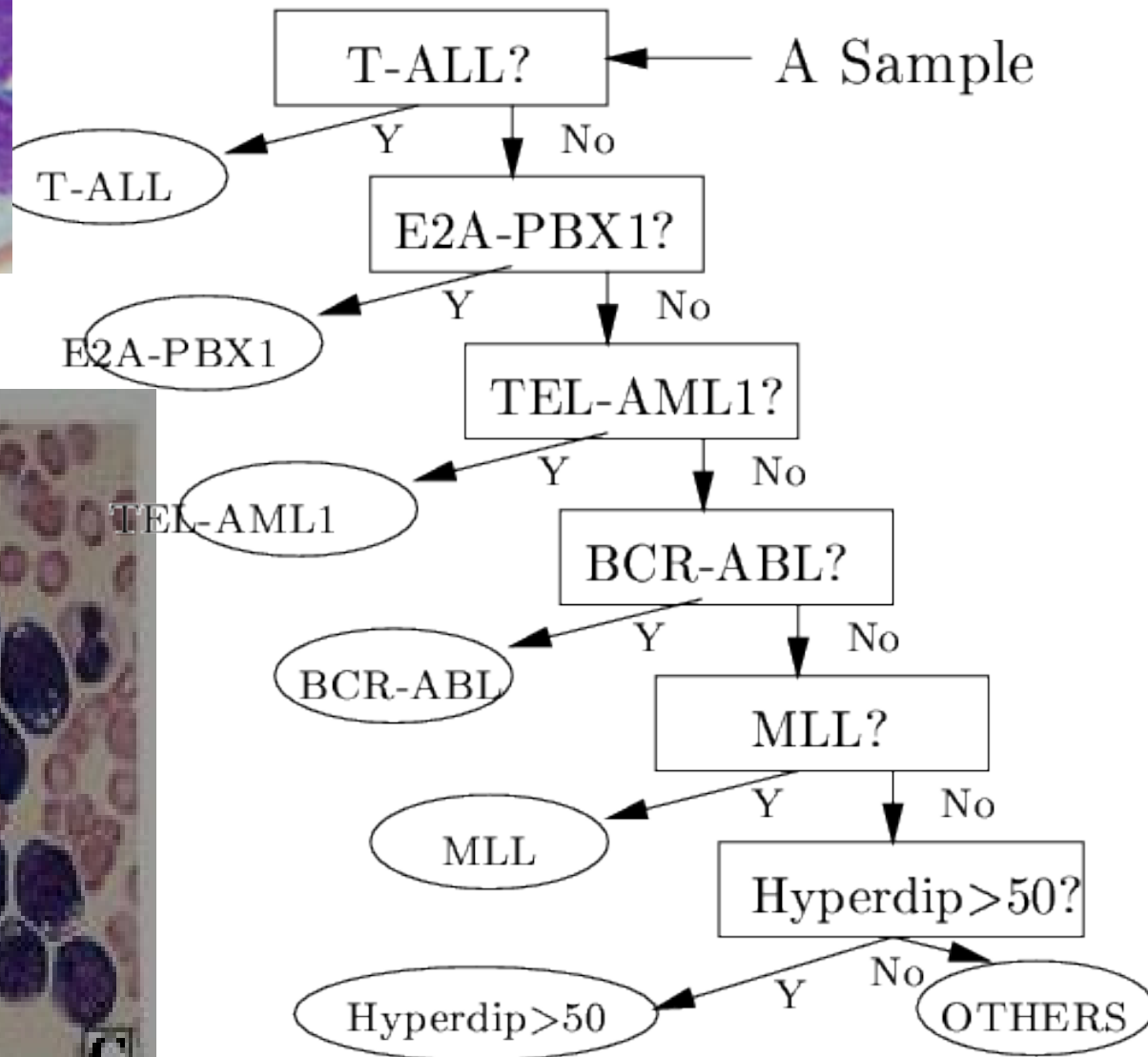
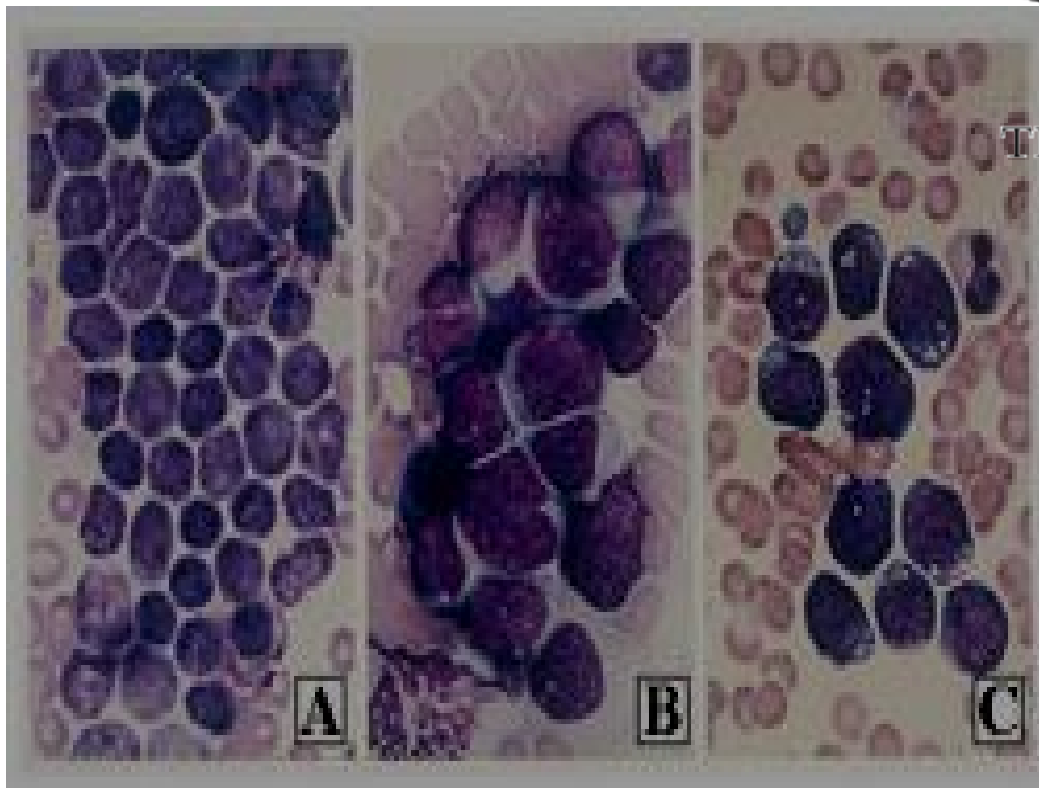
9 Motifs Discovered by MEME algo in Histone Promoter 5' Region [-250,-1] among 127 histone promoters

MOTIF NO.	MOTIF DEFINITION	TFBS AND ASSOCIATED FACTORS	TRANSFAC SITE NUMBER
1	TCTGATTGGTTA	CCAAT-box: H1TF2 (La Bella et al. 1989; Martinelli and Heintz 1994; Gallinari et al. 1989), HiNF-B (van Wijnen et al. 1988a,b), NF-Y (Mantovani 1999), HiNF-D (van Wijnen et al 1996; Grimes et al. 2003)	R00660
2	ATGCAAATGAGG	Oct-1: Octamer transcription factor 1 (OTF-1) (Fletcher et al. 1987)	R00662
3	CTATAAAAACC	TATA-box: TBP, TFIID (Nakajima et al. 1988)	R00770
4	TTTTCGCGCCCA	E2F-binding site: E2F-1 factor (Oswald et al. 1996)	R09798
5	CAATCAGGTCCG	H4TF2 binding site: H4TF2 (La Bella and Heintz 1991)	R00681
6	AACAAACACAA	AC-box: H1TF1 (La Bella et al. 1989), HiNF-A (van Wijnen et al. 1988b), HiNF-D (van Wijnen et al. 1996; Grimes et al. 2003)	R00658
7	CAGCCAATCAGA	CCAAT-box: H1TF1 (La Bella et al. 1989), HiNF-B (van Wijnen et al. 1988a,b), NF-Y (Mantovani 1999), HiNF-D (van Wijnen et al. 1996; Grimes et al. 2003), H1TF2 (La Bella et al. 1989; Martinelli and Heintz 1994; Gallinari et al. 1989)	R00659, R00660
8	CCATTGGTTAAA	CCAAT-box: H1TF2 (La Bella et al. 1989; Martinelli and Heintz 1994; Gallinari et al. 1989), HiNF-B (van Wijnen et al. 1988a,b), NF-Y (Mantovani 1999), HiNF-D (van Wijnen et al. 1996; Grimes et al. 2003)	R00660
9	CCCCGCCCCCCG	GC-box: HiNF-C (van Wijnen et al. 1989), Sp1 (Courey and Tjian 1988), Sp3 (Bimbaum et al. 1995; Hagen et al. 1994)	R00684

Diagnosis of Childhood Acute Lymphoblastic Leukemia (ALL) and Optimization of Risk-Benefit Ratio of Therapy

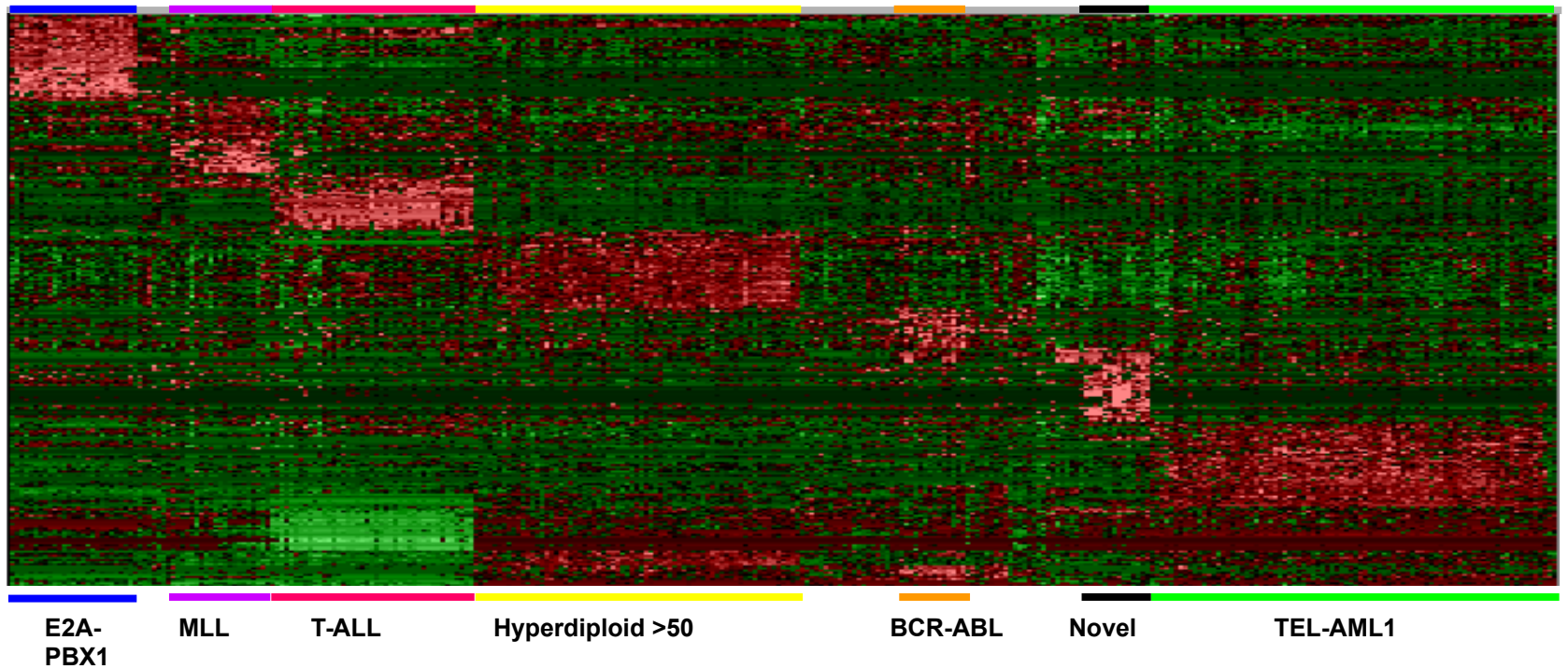


Immuno-phenotyping



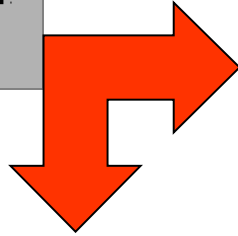


Affymetrix GeneChip Micro Array Analysis



Proteomics Data : Guilt-by-Association

Compare T with seqs of known function in a db



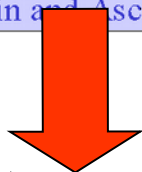
Poor Sequence Alignment

- Poor seq alignment shows few matched positions
⇒ The two proteins are not likely to be homologous

Alignment by FASTA of the sequences of amicyanin and domain 1 of ascorbate oxidase

```
Amicyanin      60      70      80      90     100
MPHNVHFVAGVLGEAALKGPMMKKEQAYSLTFTEAGTYDYHCTPHPFMRGKVVV
Ascorbate Oxidase ILQRGTPWADGTASISQCAINPGETFFYNFTVDNPGTFFYHGHLGMQRSAGLYG
                    70      80      90     100     110
```

No obvious match between
Amicyanin and Ascorbate Oxidase



Discard this function as a candidate

Good Sequence Alignment

- Good alignment usually has clusters of extensive matched positions
⇒ The two proteins are likely to be homologous

```
>gi113476732|ref|NP_108301.1| unknown protein [Mesorhizobium loti]
gi114027493|dbj|BAB53762.1| unknown protein [Mesorhizobium loti]
Length = 105

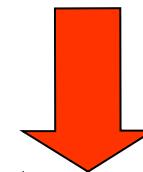
Score = 105 bits (262), Expect = 1e-22
Identities = 61/106 (57%), Positives = 73/106 (68%), Gaps = 1/106 (0%)

Query: 1  MKPGRLASIALAIIFLPMAVPAHAATIEITMENLVISPTVEVSAKVGDTIRWVVKDVFAHT 60
          MK G L  ++      MA PA AATIE+T++ LV SP  V AKVGDTI WVN DV AHT
Sbjct: 1  MKAGALIRLSWLAALALMAAPAAAATIEVTIDKLVFSPATVEAKVGDTIEWVNNDDVVAHT 60
```

good match between
Amicyanin and unknown M. loti protein

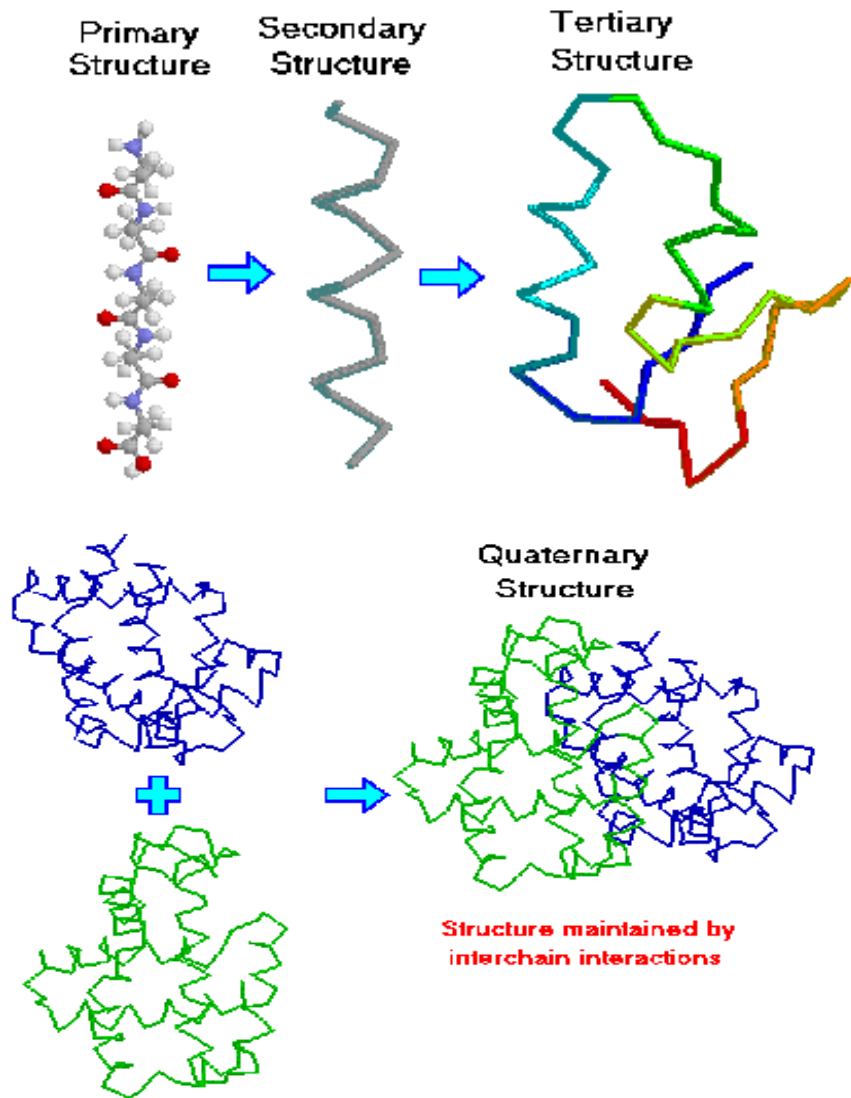


Assign to T same function as homologs



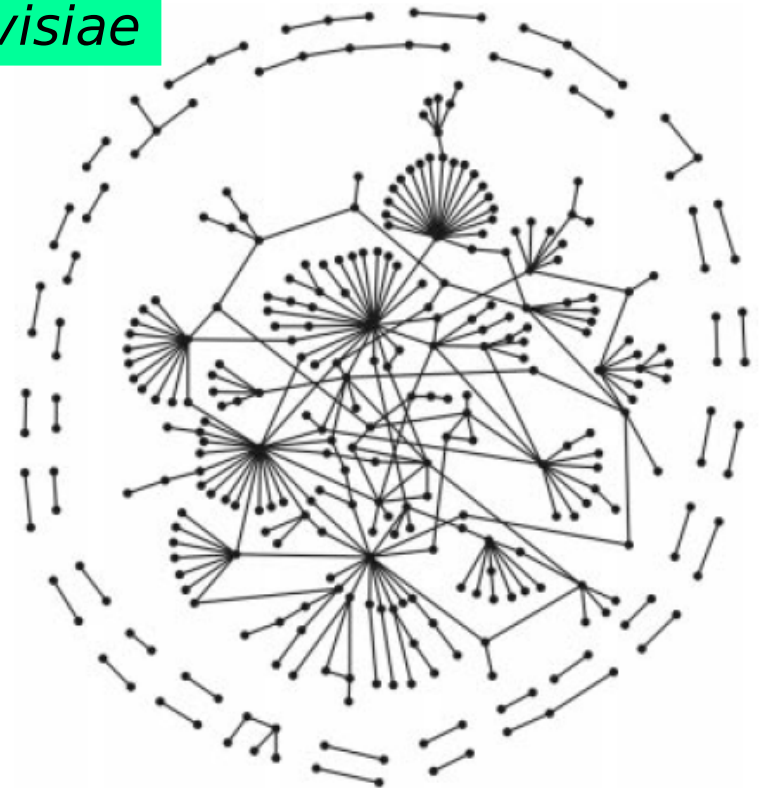
Confirm with suitable wet experiments

Proteomics Data: Subgraphs in Protein Interaction

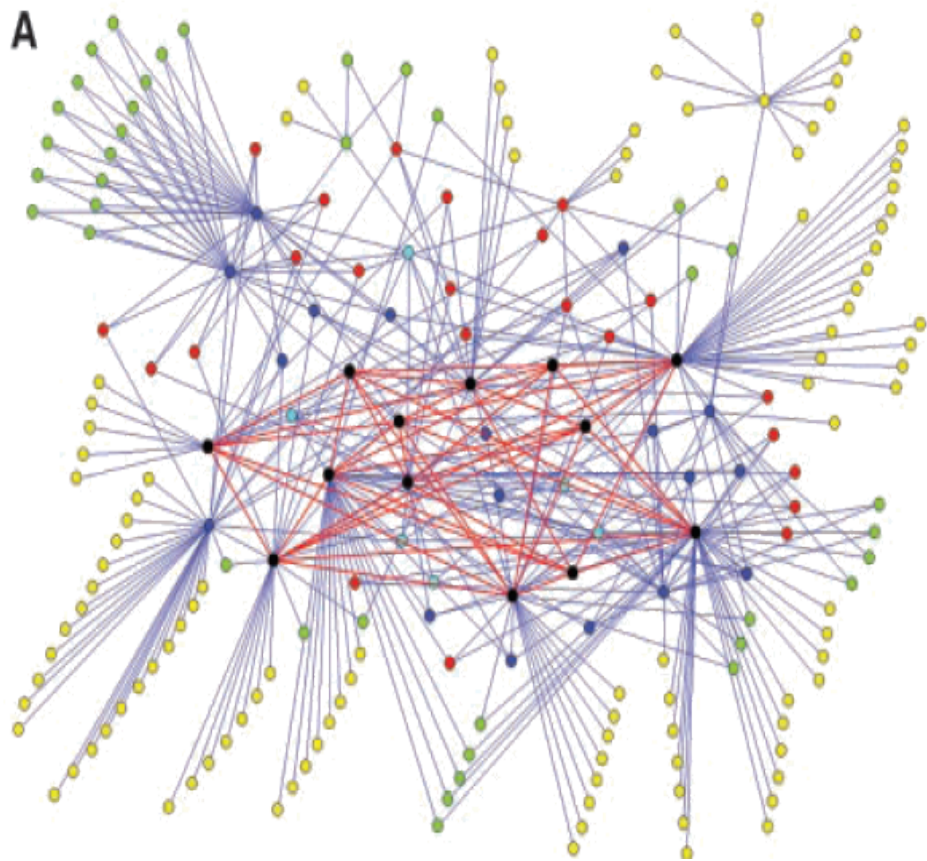


318 edges
329 nodes
In nucleus of
S. cerevisiae

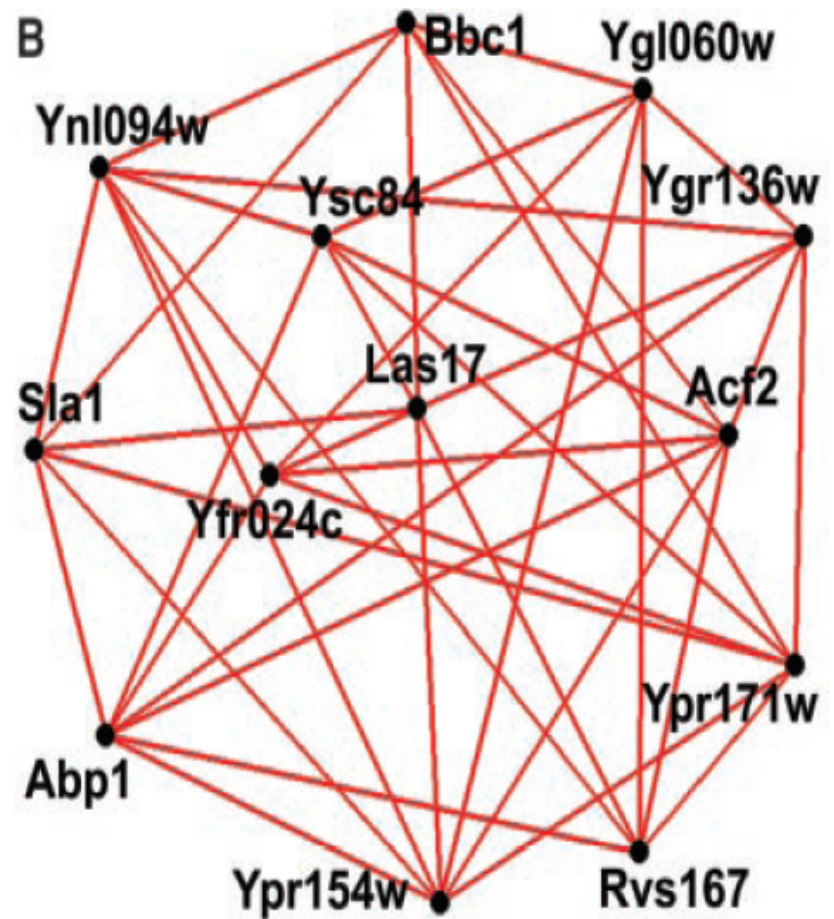
Maslov & Sneppen,
Science, v296, 2002



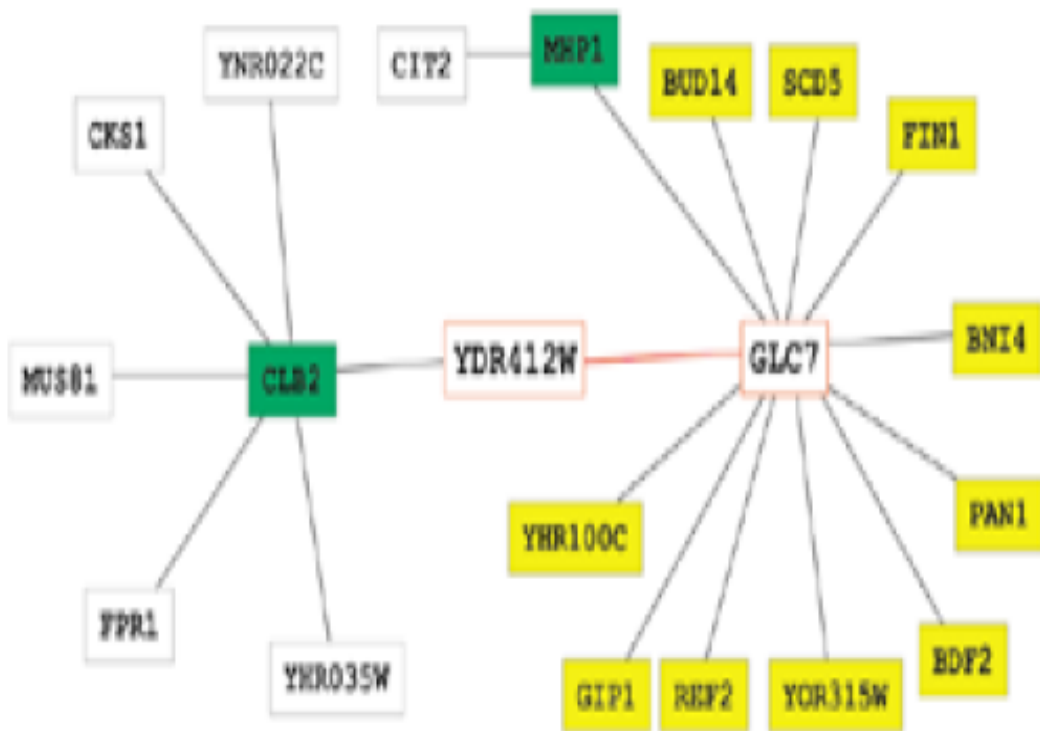
Topology of Protein Interaction Networks:
Hubs, Cores, Bipartites



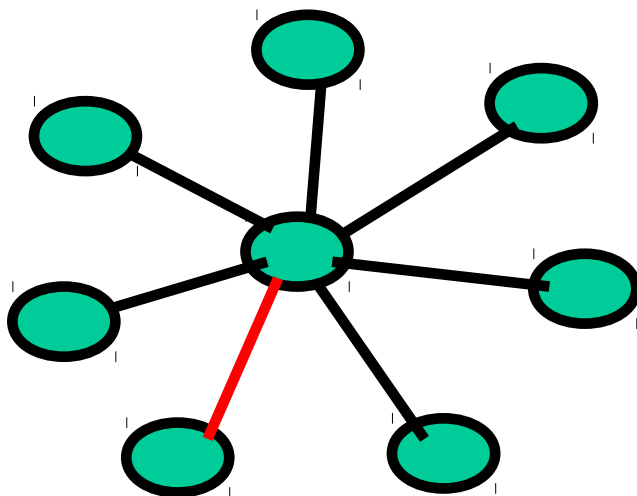
Yeast SH3 domain-domain
Interaction network:
394 edges, 206 nodes
Tong et al. *Science*, v295. 2002



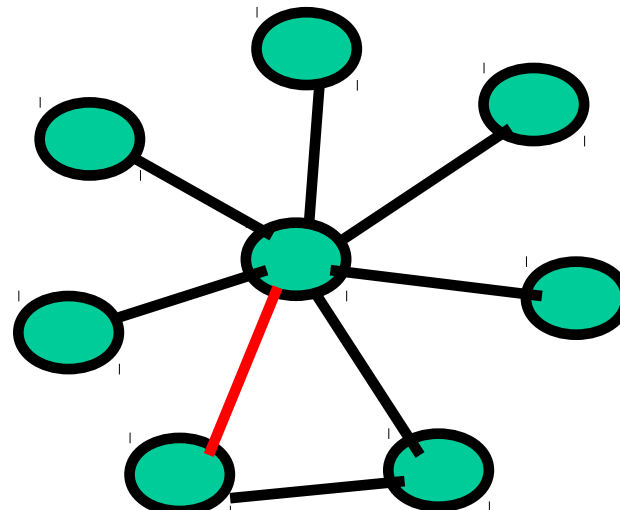
8 proteins containing SH3
5 binding at least 6 of them



Network	Nodes	Edges	N_{real}	$N_{\text{rand}} \pm \text{SD}$	Z score	N_{real}	$N_{\text{rand}} \pm \text{SD}$	Z score
Gene regulation (transcription)				Feed-forward loop			Bi-fan	
<i>E. coli</i>	424	519	40	7 ± 3	10	203	47 ± 12	13
<i>S. cerevisiae</i> *	685	1,052	70	11 ± 4	14	1812	300 ± 40	41



a

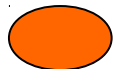
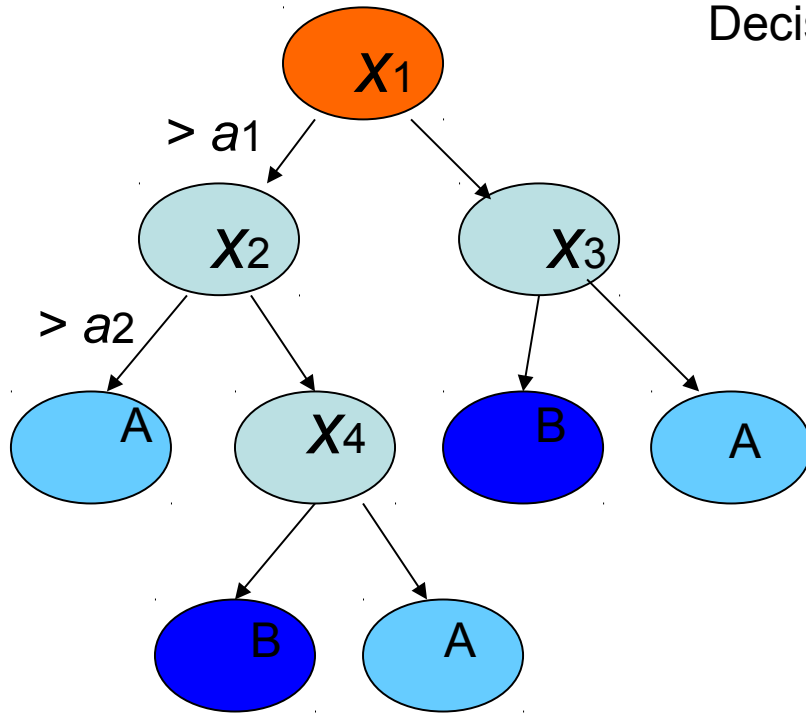


b

Configuration **a** is less likely than **b** in protein interaction networks → Graph/Network Mining

Prognosis based on Gene Expression Profiling

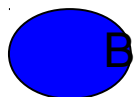
Decision Tree Based In-Silico Cancer Diagnosis



Root node



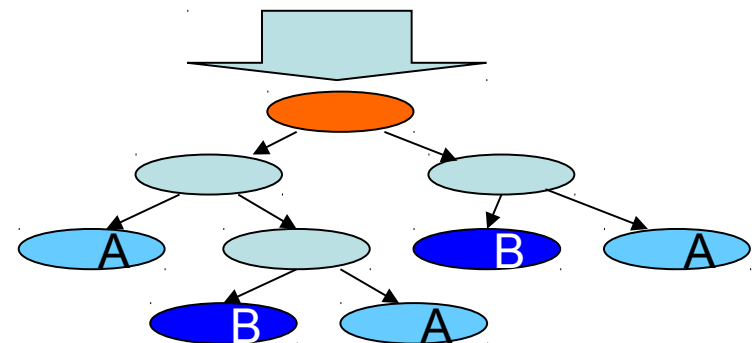
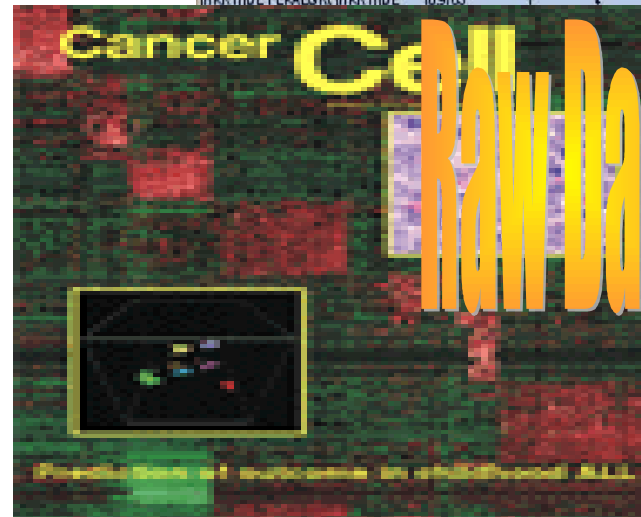
Internal nodes



Leaf nodes

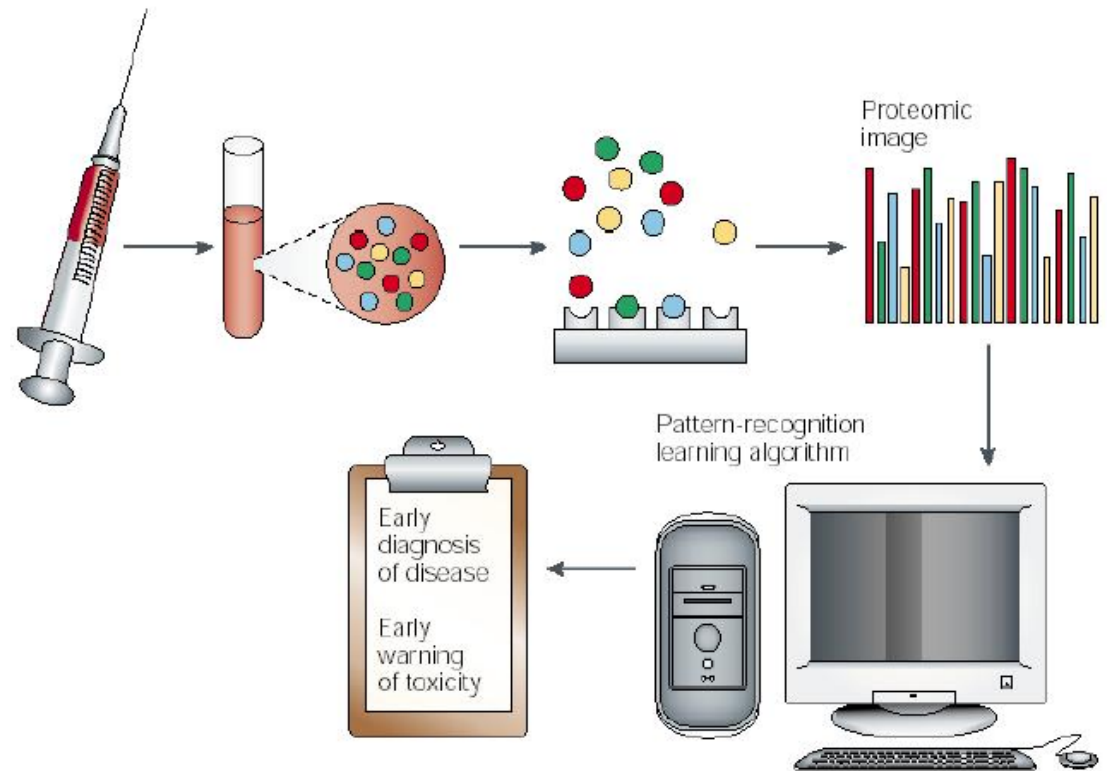


GKIQVVS	[EHM]ILGKIQVVS	3.7452	1065	3	82	3.823908	7.50615	0.0259084	0.2935175	1.8660625
P[MD]P[P]	P[LD]P[PL	4.3248	12	13	6	6.20412	5.602059	0.0384615	0.4540838	1.1737896
AAS[F]	P[MEP]IV	3.1725	37	100	65	5.602059	4.90309	0.0175676	0.1845499	1.1114275
O[EH]EG	[QD]G[M]O	3.6243	38	11	10	5.602059	6.20412	0.0239234	0.2824445	0.9382602
[EH]Q[LP	[M]EH[K]F	2.3216	126	317	392	5.124938	4.60206	0.0096142	0.095463	0.8223861
GV[F]S	P[EH]P	1.8912	795	92	534	3.60206	4.90309	0.0073011	0.0620967	0.562639
O[Q]Q	[ST]O[EH]A	1.655	130	388	313	4.60206	4.60206	0.0062054	0.0571152	0.4734859
L[F]Q[VLK	L[F]Q[MLK	2.8855	23	12	4	5.90309	6.20412	0.0144928	0.1754668	0.3509336
[MT]M[K	A[EH]M[F]	2.4993	3	150	5	6.602059	4.60206	0.0111111	0.1244902	0.2890573
P[LM]F[Q]P	P[LM]F[Q]P	2.346	25	20	5	5.590907	5.726998	0.01	0.1127191	0.2518648
IKR[TD]E	FEKAEGR[IKR]TD	10.9783	1	4	4	22.11751	10.864369	1	32.981779	65.963588



Discovery of Diagnostic Biomarkers for Ovarian Cancer

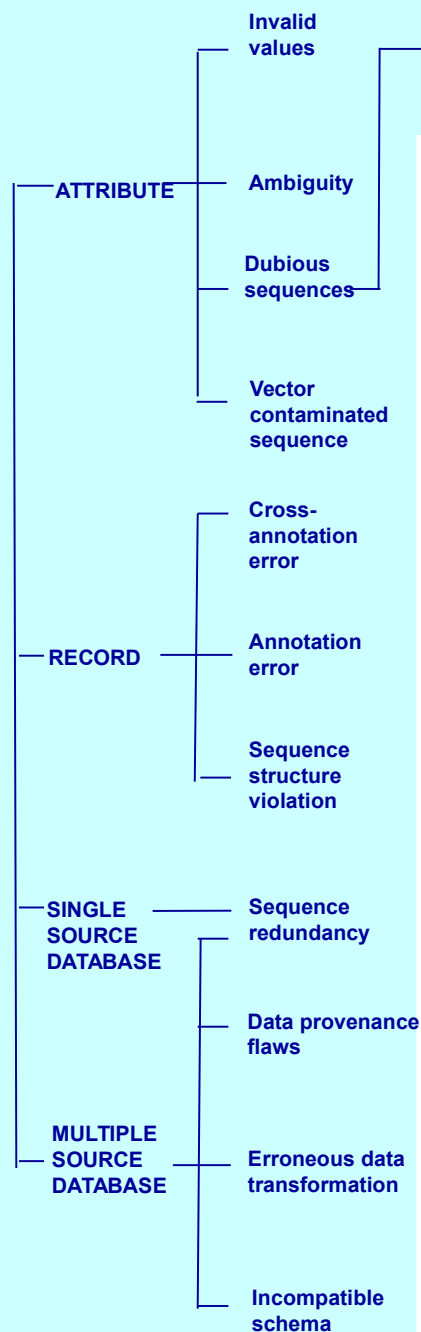
- Motivation: cure rate ~ 95% if correct diagnosis at early stage
- Proteomic profiling data obtained from patients' serum samples
- The first data set by Petricoin et al was published in *Lancet*, 2002
- Data set of June-2002.
- 253 samples: 91 controls and 162 patients suffering from the disease; 15154 features (proteins, peptides, precisely, mass/charge identities)



Methods	CS4	C4.5		
		Single	Bagging	Boosting
Errors	0 (0:0)	10 (4:6)	7 (3:4)	10 (4:6)

SVM: 0 errors; Naïve Bayes: 19 errors; k-NN: 15 errors.

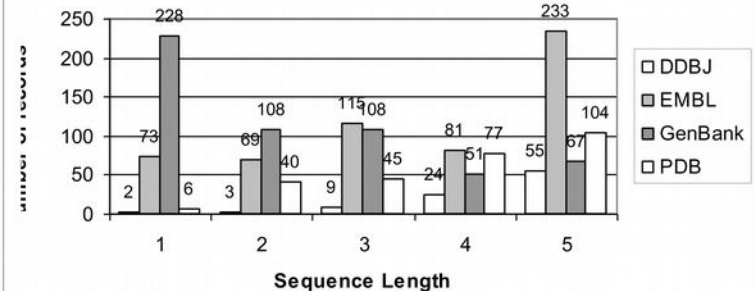
Mining Errors from Bio Databases

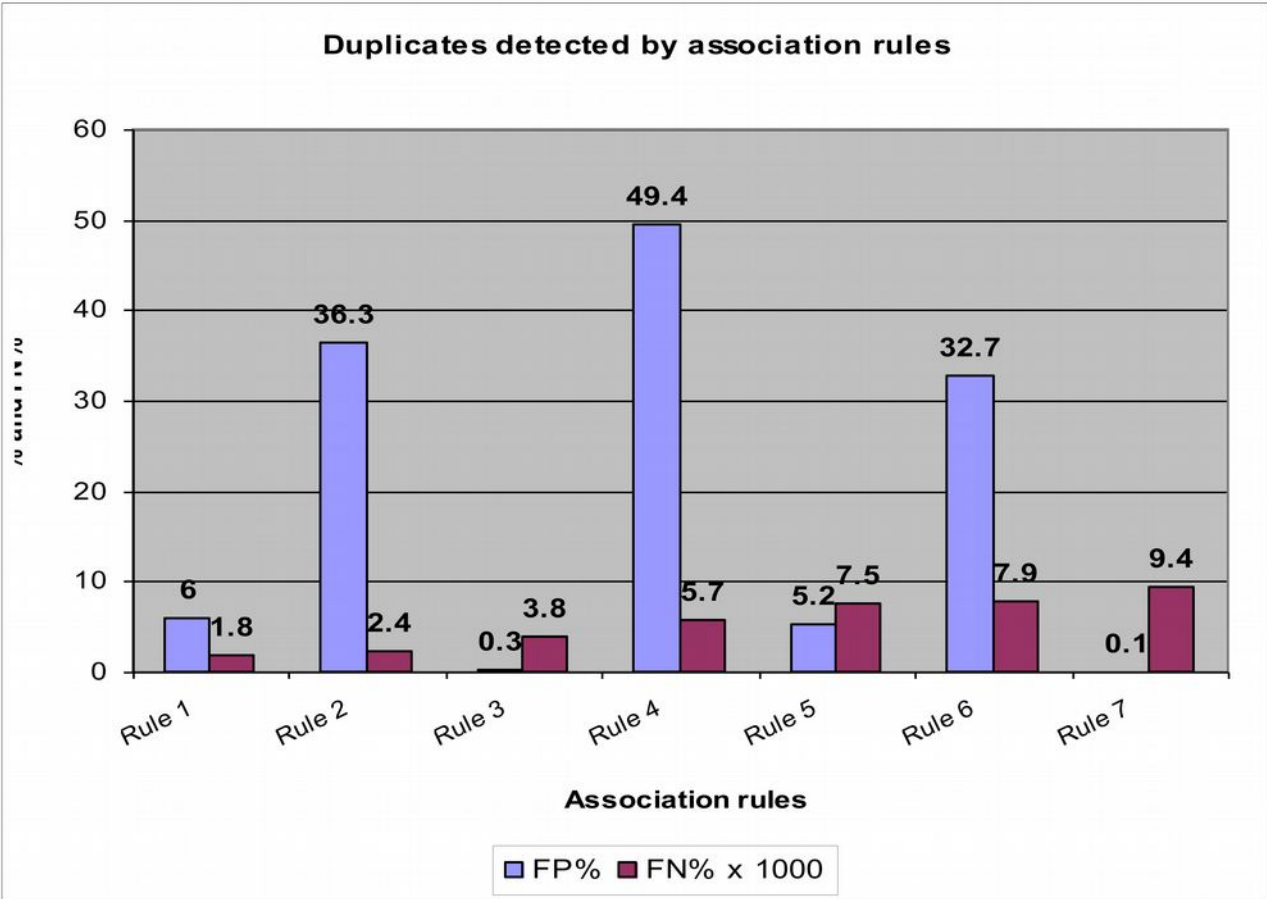


Example Meaningless Seqs

- Among the 5,146,255 protein records queried using Entrez to the major protein or translated nucleotide databases , **3,327** protein sequences are shorter than four residues (as of Sep, 2004).
- In Nov 2004, the total number of undersized protein sequences increases to **3,350**.
- Among 43,026,887 nucleotide records queried using Entrez to major nucleotide databases, **1,448** records contain sequences shorter than six bases (as of Sep, 2004).
- In Nov 2004, the total number of undersized nucleotide sequences increases to **1,711**.

Undersized nucleotide sequences in major databases



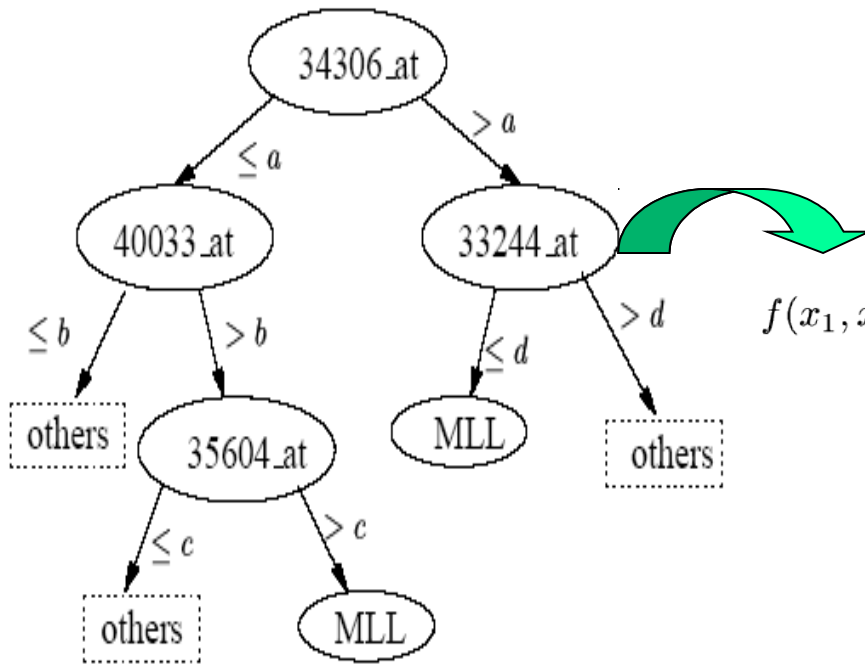


Rule 1	$S(\text{Seq})=1 \wedge N(\text{Seq Length})=1 \wedge M(\text{PDB})=0$ (99.7%)
Rule 2	$S(\text{Seq})=1 \wedge M(\text{PDB})=0 \wedge M(\text{Species})=1$ (97.1%)
Rule 3	$S(\text{Seq})=1 \wedge N(\text{Seq Length})=1 \wedge M(\text{Species})=1 \wedge M(\text{PDB})=0$ (96.8%)
Rule 4	$S(\text{Seq})=1 \wedge M(\text{PDB})=0 \wedge M(\text{DB})=0$ (93.1%)
Rule 5	$S(\text{Seq})=1 \wedge M(\text{Seq Length})=1 \wedge M(\text{PDB})=0 \wedge M(\text{DB})=0$ (92.8%)
Rule 6	$S(\text{Seq})=1 \wedge M(\text{Species})=1 \wedge M(\text{PDB})=0 \wedge M(\text{DB})=0$ (90.4%)
Rule 7	$S(\text{Seq})=1 \wedge N(\text{Seq Length})=1 \wedge M(\text{Species})=1 \wedge M(\text{PDB})=0 \wedge M(\text{DB})=0$ (90.1%)

Rule 1. Identical sequences with the same sequence length and not originated from PDB are 99.7% likely to be duplicates.

Rule 2. Identical sequences with the same sequence length and of the same species are 97.1% likely to be duplicates.

Rule 3. Identical sequences with the same sequence length, of the same species and not originated from PDB are 96.8% likely to be duplicates.

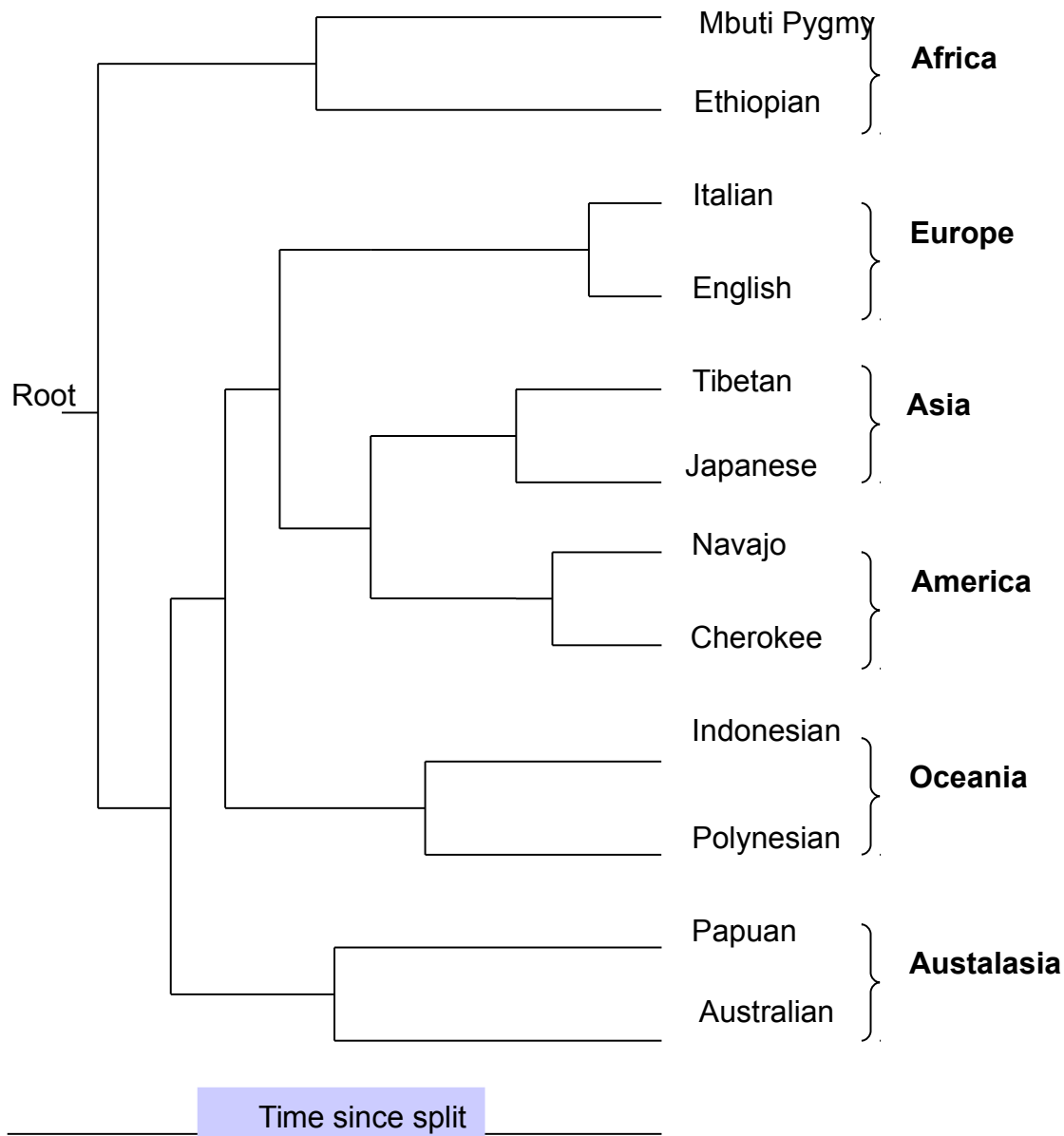


$$f(x_1, x_2, x_3, x_4) = \begin{cases} -1 & \text{if } x_1 \leq a, x_2 \leq b \\ -1 & \text{if } x_1 \leq a, x_2 > b, x_3 \leq c \\ 1 & \text{if } x_1 \leq a, x_2 > b, x_3 > c \\ 1 & \text{if } x_1 > a, x_4 \leq d \\ -1 & \text{if } x_1 > a, x_4 > d \end{cases}$$

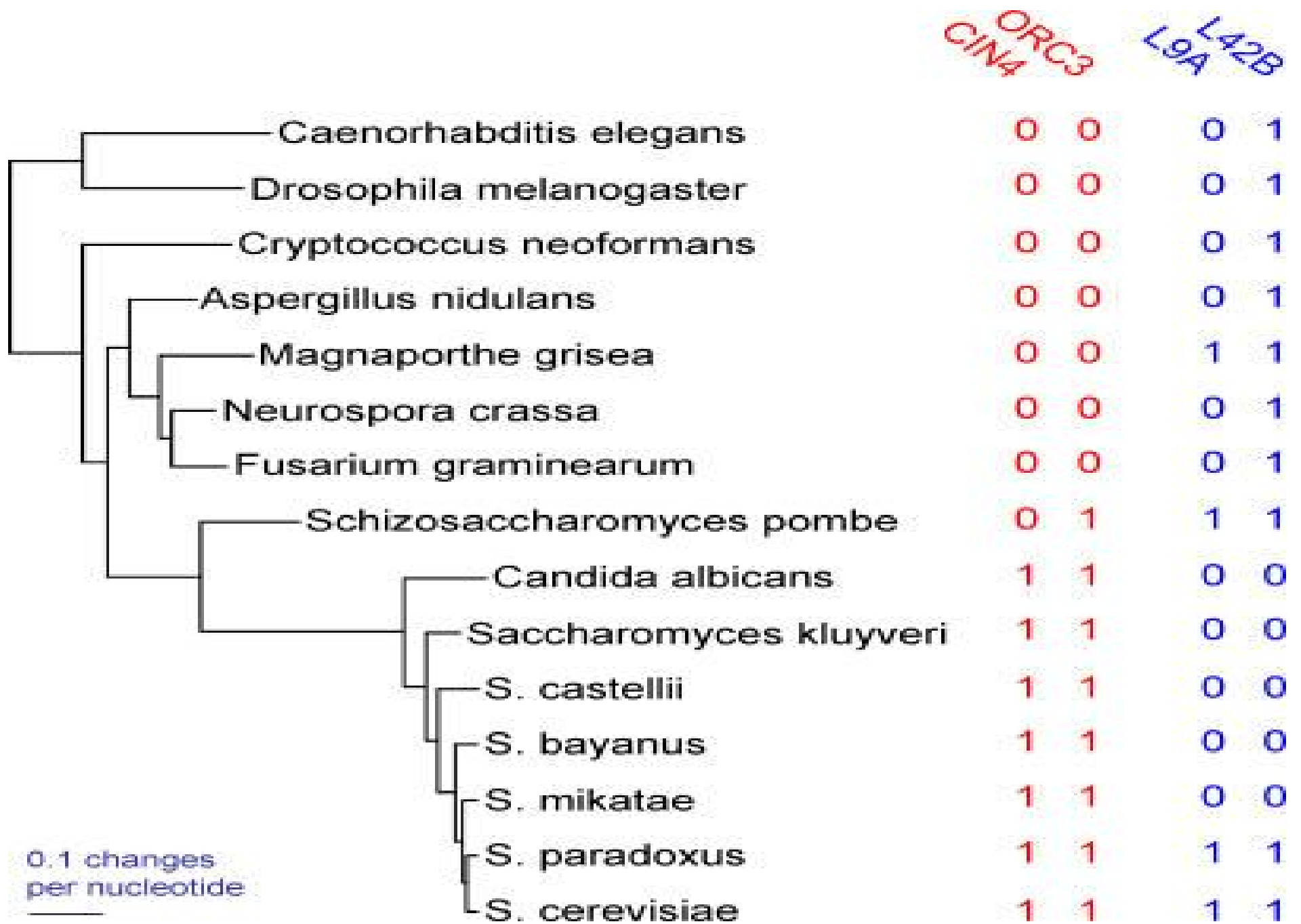
Given a test sample, at most 3 of the 4 genes' expression values are needed to make a decision!

- Yeoh et al., *Cancer Cell* 1:133-143, 2002; Differentiating MLL subtype from other subtypes of childhood leukemia
- Training data (14 MLL vs 201 others), Test data (6 MLL vs 106 others), Number of features: 12558

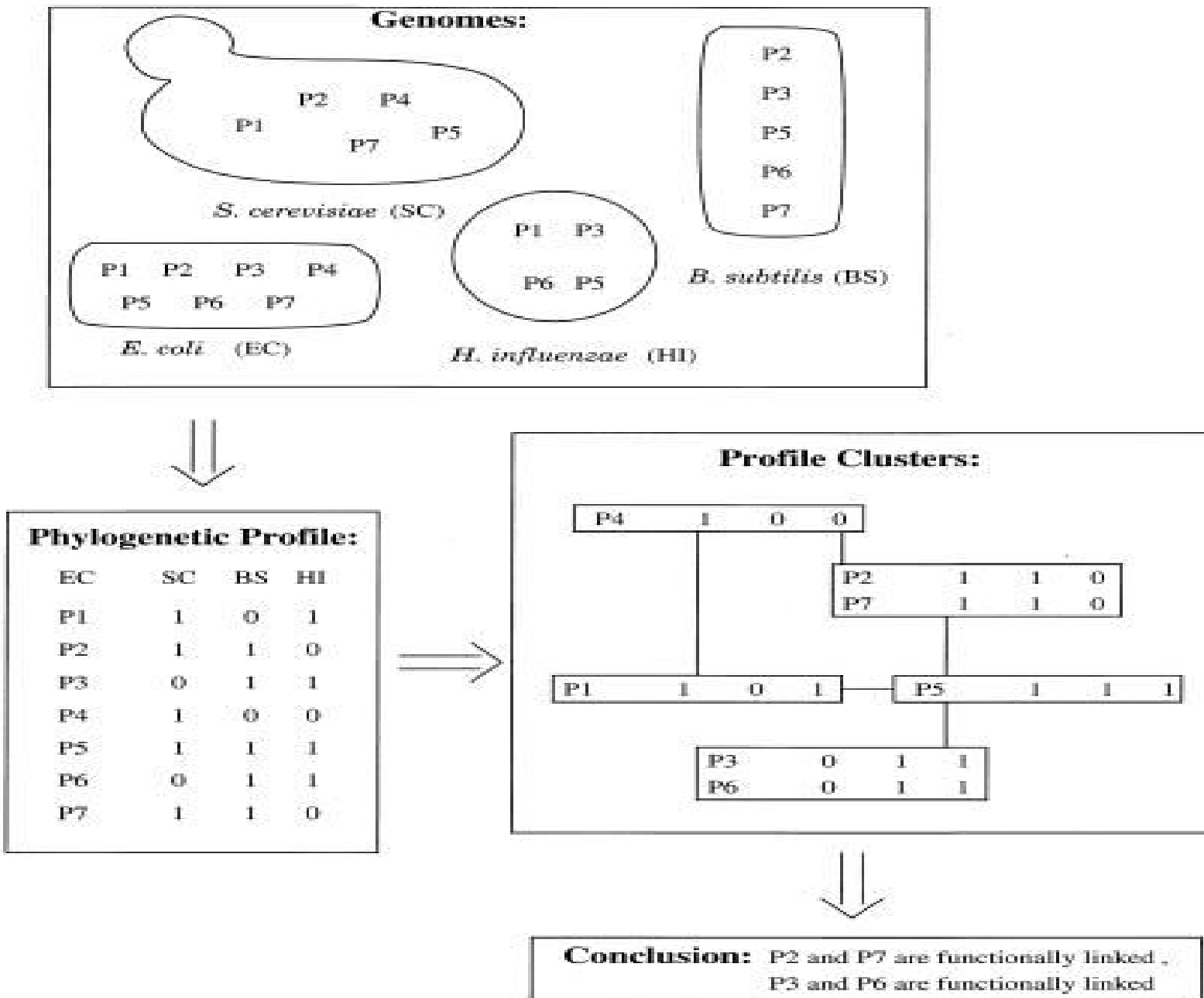
Phylogenetic tree construction



- Estimate order in which “populations” evolved
- Based on assimilated freq of many different genes
- But ...
 - is human evolution a succession of population fissions?
 - Is there such thing as a proto-Anglo-Italian population which split, never to meet again, and became inhabitants of England and Italy?



Predicting interactions using phylogenetic profile



Comparative genomics

