

Manifestation and Exploitation of Invariants in Bioinformatics

Limsoon Wong
School of Computing &
School of Medicine



Guest lecture for USP, 7/2/2007

What is invariant?



- **Suppose you have a bag of x red beans and y green beans**
- **Repeat the following:**
 - Remove 2 beans
 - If both green, discard both
 - If both red, discard one, put back one
 - If one green and one red, discard red, put back green
- **If one bean is left behind, can you predict its colour?**

Guest lecture for USP, 7/2/2007

Copyright 2007 © Limsoon Wong

Plan



- **Invariants in Evolution**
 - Finding Active Sites
- **From Invariants to Emerging Patterns**
 - Finding Key Mutation Sites
- **From Invariants to Origin of Species**
 - Where do Polynesians come from
- **From Invariants to “Guilt by Association”**
 - Predicting Protein Functions
- **Invariants in Diseases**
 - Identifying ALL subtypes

Invariants in Evolution



What is a domain

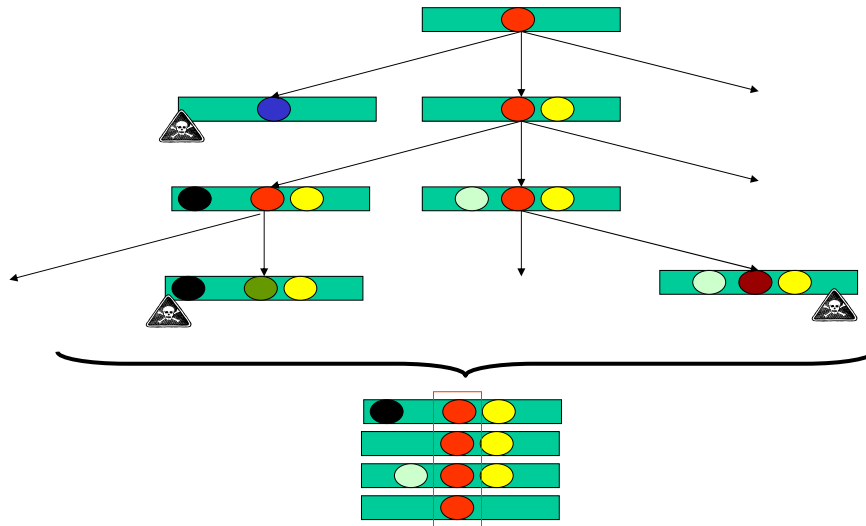
- A **domain** is a component of a protein that is self-stabilizing and folds independently of the rest of the protein chain
 - Not unique to protein products of one gene; can appear in a variety of proteins
 - Play key role in the biological function of proteins
 - Can be "swapped" by genetic engineering between one protein and another to make chimeras
- May be composed of one, more than one, or not any **structural motifs** (often corresponding to **active sites**)

Discovering Domain and Active Sites

```
>gi|475902|emb|CAA83657.1| protein-tyrosine-phosphatase alpha
MDLWFFVLLGSGLSVIGATNVITTEPPTTVPSTRIPTKAPTAAPDGGTTPRVSSSLNVSSPMTTSAPASE
PPTTTATSI SPNATTASLNASTPGTSVPTSAPVAISLPPSATPSALLTALPSTEAMTERNVSATVTTQE
TSSASHNGNSDRREDETPIIAVMVALSLLVIVFIIIVLYMLRFKYYKQAGSHSNSFRLPNGRTDDAEPQS
MPLLARSPTNRKYPPPLVVDKLEEEINRRIGDDNKLFRREFNALPACPIQATCEAASKEENKEKNRYVNI
LPYDHSRVHLTPVEGVPSHSHYINTSFINSYQEKNFIAAQGPKEETVNDVFRMIWEQNTATIVMVTNLKE
RKECKCAQYWPDQGCWTYGNIRVSVEDVTVLVDYTVRKFCIQQVGDVTNKKPQRLVTQPHFTSWPDFGVP
FTP IGMFLKFLKKVKTCPNPQYAGAI VVHCSAGVGRGTGTFIVIDAMLDMMHAERKVDVYGFVSRIRAQRQCM
VQTDQMYYVFIYQALLEHYLYGDTELEVTSLEIHLQKIYNKVPGTSSNGLEEEFKKLTISIKI QNDKMRGTGN
LPANMKKNRVLQIIPYEFNRVLIIPVKRGEENTDYVNASFIDGYRRRTPTCQPRPVQHTIEDFWRMIEWEK
SCSIVMLTELEERGQEKCAQYWP SDGSVSYGDI NVELKKEEECESYTVRDLLVTNTRENKSRQIRQFHFH
GWPEVGI PSDGKGMINI IAAVQKQQQSGNHPMHCHCSAGAGRTGTF CALSTVLERVKAEGIL DVPQTVK
SLRLQRPHMVQTLEQYEFYKVVQEYIDAFSDYANFK
```

- How do we find the domain and associated active sites in the protein above?

In the course of evolution...



Multiple Alignment of PTPs

```

gi|126467|      FHFTSWPDFGVFPFTP I GMLKFLKVKACNP--QYAGAI VVHCSAGVGRGTGTFVVIDAMLD
gi|2499753     FHFTGWPDHGVPHYATGLLSF IRRVKLSNP--PSAGP I VVHCSAGAGRTGCYIVIDIMLD
gi|462550|     YHFTQWPDHGVPEYALPVLTFVRRSSAARM--PETGPVLVHCSAGVGRGTGYIVIDSMLQ
gi|2499751     FHFTSWPDHGVPTD TDLINFRYLVRD YMKQSPPEP I LVHCSAGVGRGTGTF I AIDRLIY
gi|1709906     FQFTAWPDHGVPEHP TFLAFLRRVKTCNP--PDAGPMVVHCSAGVGRGTGCF I VIDAMLE
gi|126471|     LHFTSWPDFGVFPFTP I GMLKFLKVKTLNP--VHAGP I VVHCSAGVGRGTGTF I VIDAMMA
gi|548626|     FHFTGWPDHGVPHYATGLLSF IRRVKLSNP--PSAGP I VVHCSAGAGRTGCYIVIDIMLD
gi|131570|     FHFTGWPDHGVPHYATGLLGFVRQVKS KSP--PNAGPLVVHCSAGAGRTGCF I VIDIMLD
gi|2144715     FHFTSWPDHGVPTD TDLINFRYLVRD YMKQSPPEP I LVHCSAGVGRGTGTF I AIDRLIY
..*  ***  ***          .  *          ..*****  *****  ** ..

```

- Notice the PTPs agree with each other on some positions more than other positions
 - These positions are more imp't wrt PTPs
 - Else they wouldn't be conserved by evolution
- ⇒ They are candidate active sites

A Twist in the Tale: From Invariant to Emerging Pattern



Guest lecture for USP, 7/2/2007

Identifying Key Mutation Sites

K.L.Lim et al., *JBC*, 273:28986--28993, 1998



Sequence from a typical PTP domain D2

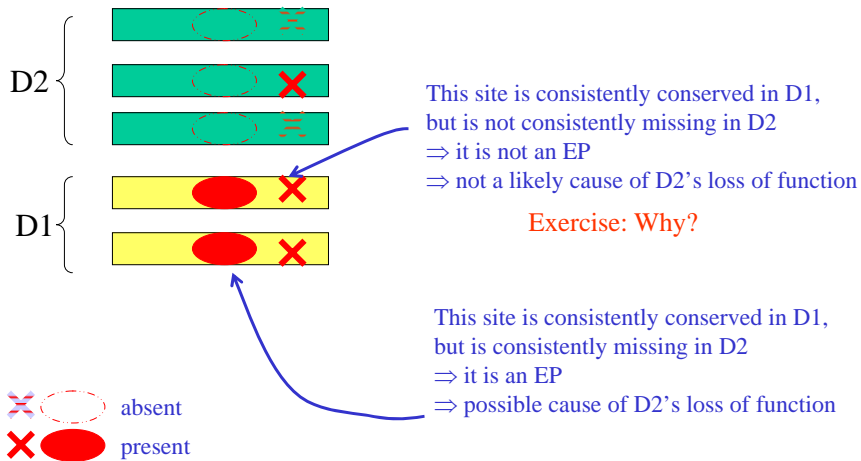
```
>g1|00000|PTPA-D2  
EEEFKILTSIKIONDKERTGMLPANEKKNVQLIIPYEFNRWIIPVKEGEMTDYVNASF  
IDQYRQDSYIASQOPLEETIEDFURHIEWRSCSIVELTELEERQQRCAQTSPSDOLV  
SYODITVELKKEEKECESTTVRDLVYNTREKESRQIRQFDFBOUPEVQIPSDGKQKDSII  
AAVQRQOQOQONEPITVBCSAGAGRTOTTFCALSTVLERVKAEGILDVFQTVKSLRLQRPK  
EIQTLKQYEFCTVWQETIDAFSDYANFK
```

- Some PTPs have 2 PTP domains
- PTP domain D1 is has much more activity than PTP domain D2
- Why? And how do you figure that out?

Guest lecture for USP, 7/2/2007

Copyright 2007 © Limsoon Wong

Invariant as Emerging Pattern



Guest lecture for USP, 7/2/2007

Copyright 2007 © Limsoon Wong

Emerging Patterns of PTP D1 vs D2



- Collect example PTP D1 sequences
- Collect example PTP D2 sequences
- Make multiple alignment A1 of PTP D1
- Make multiple alignment A2 of PTP D2
- Are there positions conserved in A1 that are violated in A2?
- These are candidate mutations that cause PTP activity to weaken
- Confirm by wet experiments

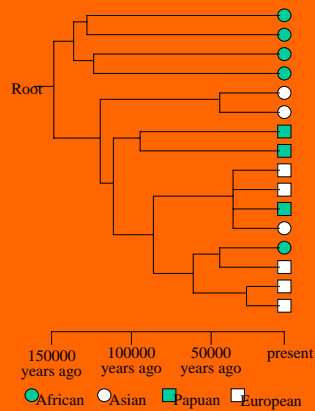
Guest lecture for USP, 7/2/2007

Copyright 2007 © Limsoon Wong

Confirmation by Mutagenesis Expt

- **What wet experiments are needed to confirm the prediction?**
 - Mutate E → D in D2 and see if there is gain in PTP activity
 - Mutate D → E in D1 and see if there is loss in PTP activity

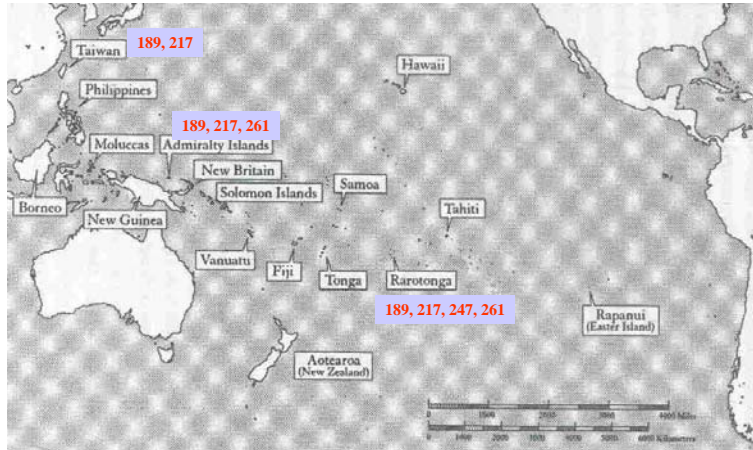
Exercise: Why do you need this 2-way expt?



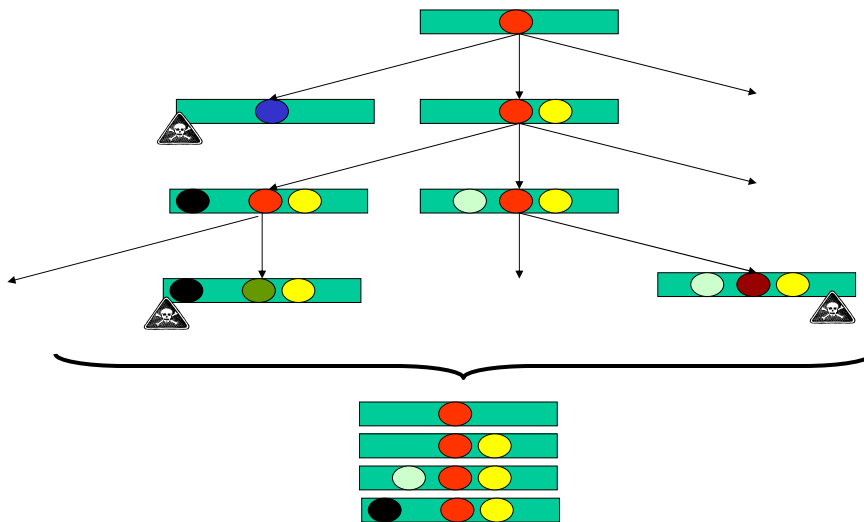
From Invariant to
Origin of Species

Origin of Polynesians

- Do they come from Asia or America?



In the course of evolution...



Origin of Polynesians

- Common mitochondrial control seq from Rarotonga have variants at positions 189, 217, 247, 261. Less common ones have 189, 217, 261
- More 189, 217 closer to Taiwan. More 189, 217, 261 closer to Rarotonga
- 247 not found in America
⇒ Polynesians came from Taiwan!
- Seq from Taiwan natives have variants 189, 217
- Taiwan seq sometimes have extra mutations not found in other parts
⇒ These are mutations that happened since Polynesians left Taiwan!
- Seq from regions in betw have variants 189, 217, 261.

From Invariant to Guilt by Association

A protein is a ...

- A protein is a large complex molecule made up of one or more chains of amino acids
- Protein performs a wide variety of activities in the cell



Function Assignment to Protein Sequence

```
SPSTNRKYPPLPVDKLEEEINRRMADDNKLFREEFNALPACPIQATCEAASKEENKEKNR
YVNILPYDHSRVHLTPVEGVPSDYINASFINGYQEKNKFIAAQGPKEETVNDFWRMWE
QNTATIVMTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD
VTNRKPQRLITQFHFTSWPDFGVPFPIGMLKFLKKVKACNPQYAGAIVVHCSAGVGRGTG
TFVVIDAMLDMMHSEKVDVYGFVSRIRAQRCQMVTDMQYVFIYQALLEHYLYGDTELE
VT
```

- How do we attempt to assign a function to a new protein sequence?

Sequence Alignment: Poor Example

- 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

```

                60      70      80      90      100
Amicyanin      MPHNVHFVAGVLGEAALKGPMKKEQAYSLTFTEAGTYDYHCTFHPFMRGKVVVE
                :..: . :. :. :
Ascorbate Oxidase ILQRGTFWADGTASISQCAINPGETFFYNFTVDNPGTFFFYHGHLGMQRSAGLYGSLI
                70      80      90      100      110      120
  
```

No obvious match between
Amicyanin and Ascorbate Oxidase

Sequence Alignment: Good Example

- 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  MKPGRLASIALAIIFLPMVPAHAATIEITMENLVISPTVEVSAKVGDTIRVWNKDVFAHT 60
          MK G L  ++      MA PA AATIE+T++ LV SP  V AKVGDITVWN DV AHT
Sbjct: 1  MKAGALIRLSWLAALALMAAPAAAATIEVTIDKLVFSPATVEAKVGDTIEWVNNDDVVAHT 60
  
```

good match between
Amicyanin and unknown *M. loti* protein

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
MFRVYVFWVGLLEAALGSPFRKGGQATSLQTEAGTDFRCTYRFFPRGKRVVY
Ascorbate Oxidase ILQKQTFWADGTASISQCAINPCEYFFNFVDPOTFFYRHLGNQRSAGLYG
                    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

```

gi|11476722|ref|NP_188331.1| unknown protein [Mesorhizobium loti]
gi|14507493|tbl|TMS576.1| 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  MDPRLAALAIIFLPMWYAKHATIEITNRELVIFSTYKAWKQVDFPRKQWFAIT 60
           M G L  ++  MA FA AATIE++ LP DP  Y AKKQDTI  FVN DP AV ART
Sbjct: 1  MDAALHLEFLALALGAPAAAKATITVTLGCEATYAKVQVDFPRKQWFAIT 60
    
```

good match between
Amicyanin and unknown M. loti protein

Assign to *T* same function as homologs

Confirm with suitable wet experiments

Homologs obtained by BLAST

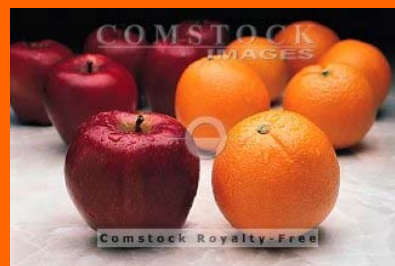
Sequences producing significant alignments:	Score (bits)	E Value
gi 14193729 cb AAK56109.1 AF332081.1 protein tyrosin phosph...	62.1	e-177
gi 126467 sp P18433 PTRA_HUMAN Protein-tyrosine phosphatase...	62.1	e-177
gi 4506303 ref INP_002827.1 protein tyrosine phosphatase, r...	62.1	e-176
gi 227294 prf I1701300A protein Tyr phosphatase	62.0	e-176
gi 18450369 ref INP_543030.1 protein tyrosine phosphatase, ...	62.1	e-176
gi 32067 emb CAA37447.1 tyrosine phosphatase precursor [Ho...	61.9	e-176
gi 285113 pir JJC1285 protein-tyrosine-phosphatase (EC 3.1....	61.9	e-176
gi 6981446 ref INP_036895.1 protein tyrosine phosphatase, r...	61.4	e-176
gi 2098414 pdb 1YFO1A Chain A, Receptor Protein Tyrosine Ph...	61.5	e-174
gi 32313 emb CAA38662.1 protein-tyrosine phosphatase [Homo...	61.1	e-174
gi 450583 cb AAB04150.1 protein tyrosine phosphatase >gi 4...	60.5	e-172
gi 6679557 ref INP_033006.1 protein tyrosine phosphatase, r...	60.1	e-172
gi 483922 cb AAA17990.1 protein tyrosine phosphatase alpha	59.9	e-170

- Thus our example sequence could be a protein tyrosine phosphatase α (PTP α)

What if there is no sequence homology?

Guilt by association of other invariants of evolution!

Guilt by Association of
Other Invariants of
Evolution!



Dissimilarities as Invariant!

	orange₁	banana₁	...
apple₁	Color = red vs orange Skin = smooth vs rough Size = small vs small Shape = round vs round	Color = red vs yellow Skin = smooth vs smooth Size = small vs small Shape = round vs oblong	...
apple₂	Color = red vs orange Skin = smooth vs rough Size = small vs small Shape = round vs round	Color = red vs yellow Skin = smooth vs smooth Size = small vs small Shape = round vs oblong	...
orange₂	Color = orange vs orange Skin = rough vs rough Size = small vs small Shape = round vs round	Color = orange vs yellow Skin = rough vs smooth Size = small vs small Shape = round vs oblong	..
...

SVM-Pairwise Framework

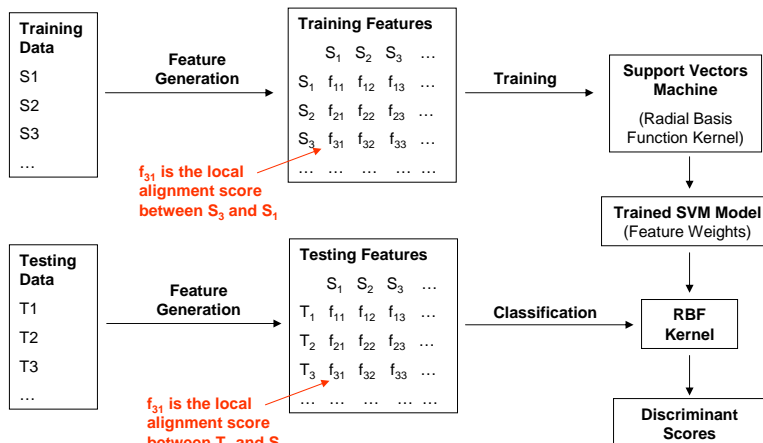
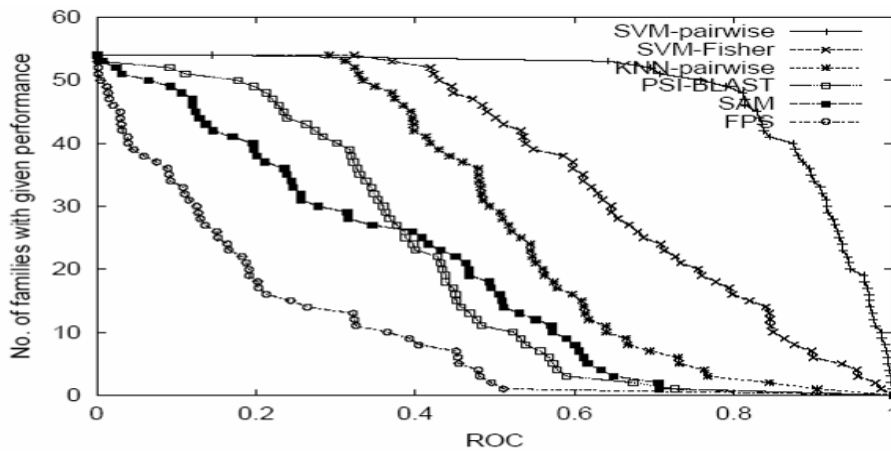


Image credit: Kenny Chua

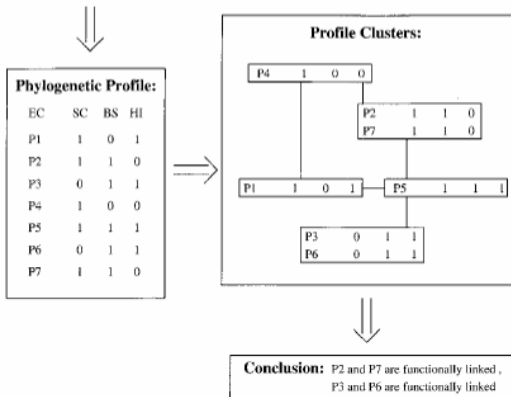
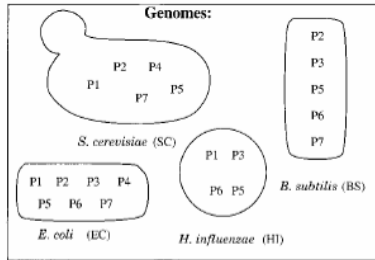
Performance of SVM-Pairwise



- **ROC: The area under the curve derived from plotting true positives as a function of false positives for various thresholds**

Phylogenetic Profiles as Invariant

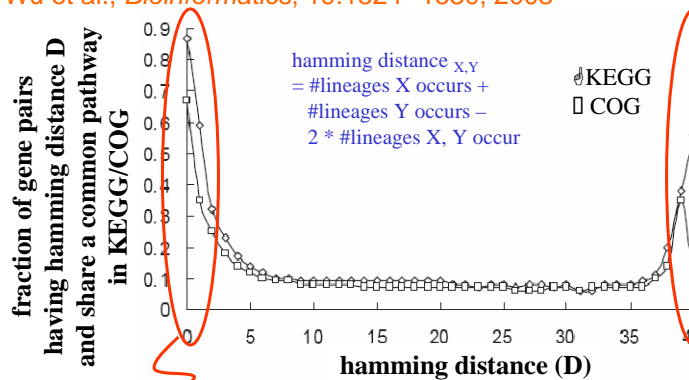
- **A protein is not alone when performing its biological function**
- ⇒ **Gene (and hence proteins) with identical patterns of occurrence across phyla tend to function together**



Phylogenetic Profiling: How It Works

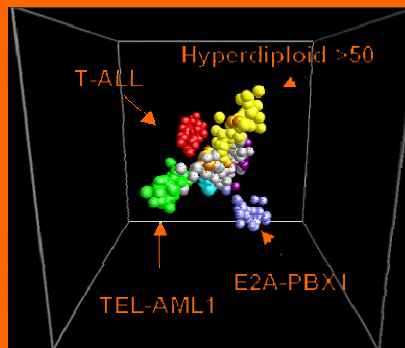
Phylogenetic Profiling: Evidence

Wu et al., *Bioinformatics*, 19:1524--1530, 2003



- Proteins having low hamming distance (thus highly similar phylogenetic profiles) tend to share common pathways
- Exercise: Why do proteins having high hamming distance also have this behaviour?

Invariants in Diseases

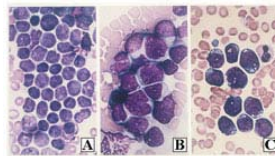


Guest lecture for USP, 7/2/2007

Childhood Acute Lymphoblastic Leukemia



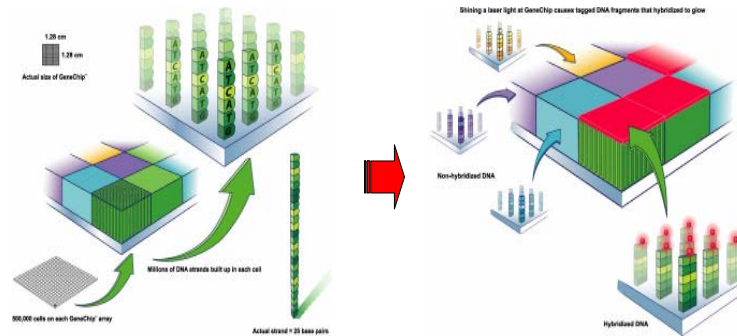
- Major subtypes: T-ALL, E2A-PBX, TEL-AML, BCR-ABL, MLL genome rearrangements, Hyperdiploid>50
- Diff subtypes respond differently to same Tx
- Over-intensive Tx
 - Development of secondary cancers
 - Reduction of IQ
- Under-intensiveTx
 - Relapse
- The subtypes look similar
- Conventional diagnosis
 - Immunophenotyping
 - Cytogenetics
 - Molecular diagnostics



Guest lecture for USP, 7/2/2007

Copyright 2007 © Limsoon Wong

Massive Gene Expression Profiling

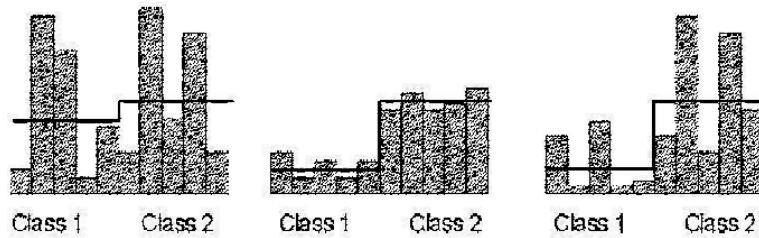


Subtype Diagnosis by Gene Expression

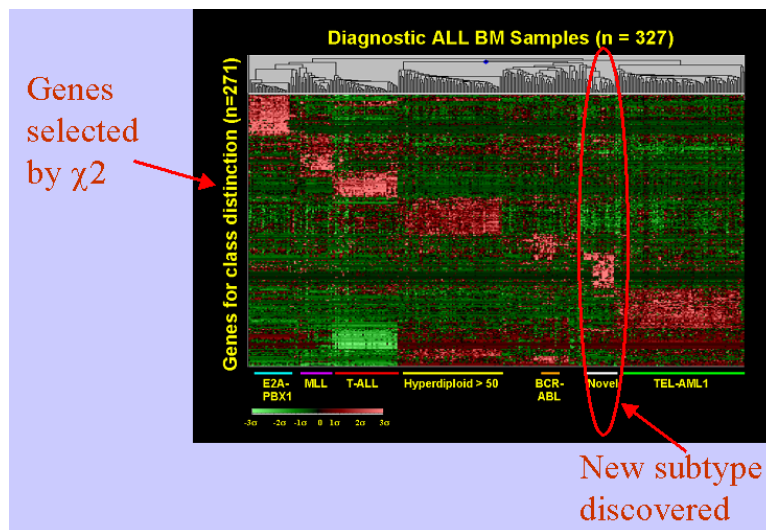
- Gene expression data collection
- Gene selection by χ^2
- Classifier training
- Apply classifier for diagnosis of future cases

Signal Selection Basic Idea

- Choose a signal w/ low intra-class distance
 - Choose a signal w/ high inter-class distance
- ⇒ Invariants which are emerging patterns



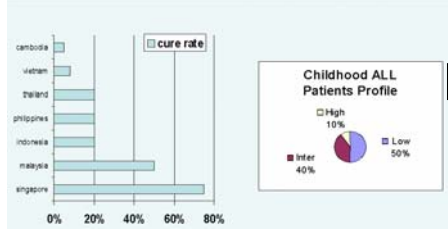
Clustering by Gene Expression Profiles



Impact

Childhood ALL in ASEAN Countries

(2000 new cases per year)



Conventional Tx:

- intermediate intensity to everyone
- ⇒ 10% suffers relapse
- ⇒ 50% suffers side effects
- ⇒ costs US\$150m/yr

Our optimized Tx:

- high intensity to 10%
- intermediate intensity to 40%
- low intensity to 50%
- costs US\$100m/yr

- High cure rate of 80%
- Less relapse
- Less side effects
- Save US\$51.6m/yr

What have we learned?

What have we learned?

- **Paradigms**
 - Invariants
 - Emerging patterns
 - “Guilt by association”
- **Techniques**
 - Sequence comparison
 - Multiple alignment
 - Machine learning
 - Signal processing
- **Applications**
 - Active sites and key mutations
 - Origin of species
 - Protein functions
 - Disease diagnosis
- **Miscellaneous**
 - Microarrays
 - Economic of bioinformatics

Suggested Readings

- Limsoon Wong, *The Practical Bioinformatician*, World Scientific, 2004. Chapters 1, 3, 4, 14.
- K.L.Lim et al. “Interconversion of kinetic identities of the tandem catalytic domains of receptor-like protein tyrosine phosphatase PTP-alpha by two point mutations is synergist and substrate dependent”, *JBC*, 273:28986--28993, 1998
- J. Wu et al. “Identification of functional links between genes using phylogenetic profiles”, *Bioinformatics*, 19:1524--1530, 2003
- T. Jaakkola, M. Diekhans, and D. Haussler. A discriminative framework for detecting remote homologies. *JCB*, 7(1-2):95—11, 2000
- B. Sykes. *The seven daughters of Eve*, Gorgi Books, 2002