# Molecular Phylogenetic Analysis
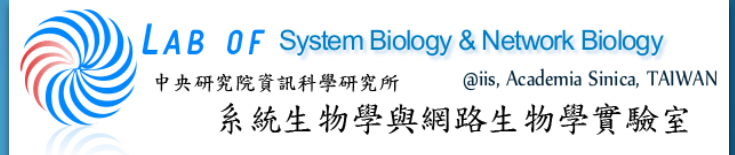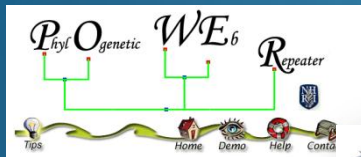
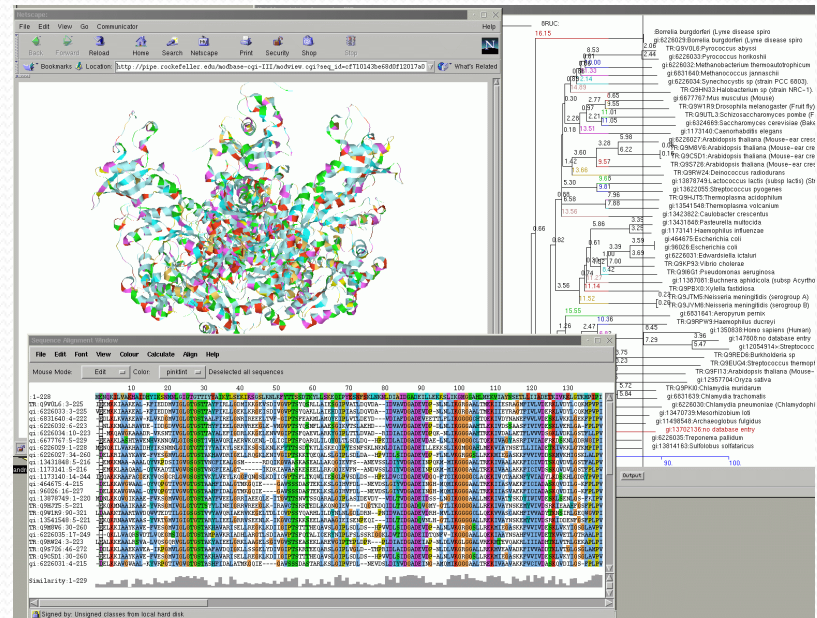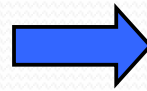林仲彥

*cylin@iis.sinica.edu.tw*
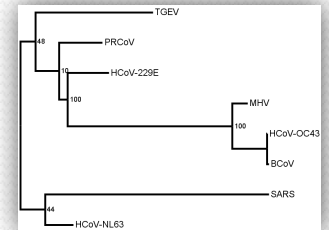
*Dec 4, 2009*

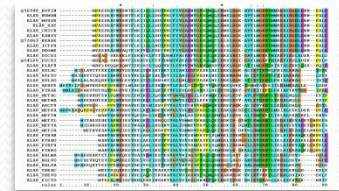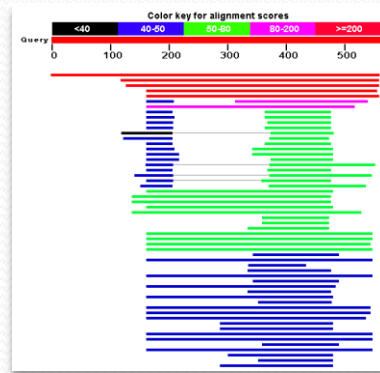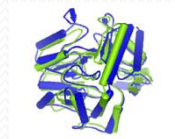*http://eln.iis.sinica.edu.tw*

# *Coding Characters and Defining Homology*



*Classical phylogenetic analysis by Morphology*

*Molecular phylogenetic analysis By Bio-Molecules*

# *Steps of Phylogenetic Analysis*



New Sequence

Homology Search

Alignment

Phylogenetic Tree

# *Elements in a Phylogenetic Tree*

- The tree is composed of nodes connected by branches.

➢ **distance scale :** scale which represents the number of differences between sequences (e.g. 0.1 means 10 % differences between two sequences)

➢ **root :** is the common ancestor of all taxa.



➢ **branch length :** often represents the number of changes that have occurred in that branch.

➢ **node :** a node represents a taxonomic unit.
   ➢ Internal nodes
   ➢ External nodes

➢ **branch (edge):** defines the relationship between the taxa.

# *Trees Only Represent The Order Of Branching*
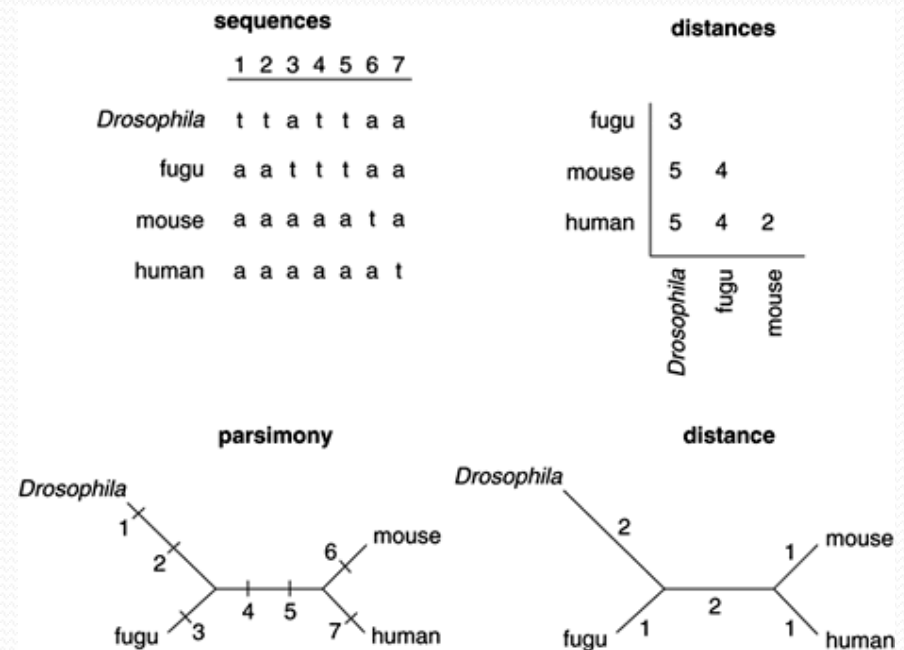
- Same topology in a different style
  - Both trees have identical topologies, with some of the internal nodes rotated.

# *The Ways to Construct the tree*

- Distance-matrix methods
  - Neighbor-joining
  - Fitch-Margoliash method
  - Using outgroups
- Maximum parsimony
  - Branch and bound
  - MALIGN and POY
- Maximum likelihood (Statistics Based)
- Bayesian inference (Probability Based)

# Phylogeny Packages

*http://evolution.genetics.washington.edu/phylip/software.html*

# *Phylip*

## ... by type of data

- DNA sequences
- Protein sequences
- Restriction sites
- Distance matrices
- Gene frequencies
- Quantitative characters
- Discrete characters
- tree plotting, consensus trees, tree distances and tree manipulation

## ... by type of algorithm

- Heuristic tree search
- Branch-and-bound tree search
- Interactive tree manipulation
- Plotting trees, consenus trees, tree distances
- Converting data, making distances or bootstrap replicat

### DNA and RNA sequence data

**DNAPARS**. Estimates phylogenies by the parsimony method using nucleic acid sequences. Allows use the full IUB ambiguity codes, and estimates ancestral nucleotide states. Gaps treated as a fifth nucleotide state. It can also fo transversion parsimony. Can cope with multifurcations, reconstruct ancestral states, use 0/1 character weights, and infer branch lengths.

**DNAMOVE**. Interactive construction of phylogenies from nucleic acid sequences, with their evaluation by parsimony and compatibility and the display of reconstructed ancestral bases. This can be used to find parsimony or compatibility estimates by hand.

**DNAPENNY**. Finds all most parsimonious phylogenies for nucleic acid sequences by branch-and-bound search. This may not be practical (depending on the data) for more than 10 or 11 species.

**DNACOMP**. Estimates phylogenies from nucleic acid sequence data using the compatibility criterion, which searches for the largest number of sites which could have all states (nucleotides) uniquely evolved on the same tree. Compatibility is particularly appropriate when sites vary greatly in their rates of evolution, but we do not know in advance which are the less reliable ones.

### Heuristic search for best tree

**PROTPARS**. Estimates phylogenies from protein sequences (input using the standard one-letter code for amino acids) using the parsimony method, in a variant which counts only those nucleotide changes that change the amino acid, on the assumption that silent changes are more easily accomplished.

**DNAPARS**. Estimates phylogenies by the parsimony method using nucleic acid sequences. Allows use the full IUB ambiguity codes, and estimates ancestral nucleotide states. Gaps treated as a fifth nucleotide state. It can also fo transversion parsimony. Can cope with multifurcations, reconstruct ancestral states, use 0/1 character weights, and infer branch lengths.

**DNACOMP**. Estimates phylogenies from nucleic acid sequence data using the compatibility criterion, which searches for the largest number of sites which could have all states (nucleotides) uniquely evolved on the same tree. Compatibility is particularly appropriate when sites vary greatly in their rates of evolution, but we do not know in advance which are the less reliable ones.

**DNAML**. Estimates phylogenies from nucleotide sequences by maximum likelihood. The model employed allows for unequal expected frequencies of the four nucleotides, for unequal rates of transitions and transversions, and for different (prespecified) rates of change in different categories of sites, and also use of a Hidden Markov model of rates, with the program inferring which sites have which rates. This also allows gamma-distribution and gamma-plus-

# *Interactive Interface for Phylip*

```
Nucleic acid sequence Maximum Likelihood method, version 3.6

Settings for this run:
  U                Search for best tree?  Yes
  T        Transition/transversion ratio:  2.0000
  F        Use empirical base frequencies?  Yes
  C                One category of sites?  Yes
  R        Rate variation among sites?  constant rate
  W                Sites weighted?  No
  S        Speedier but rougher analysis?  Yes
  G                Global rearrangements?  No
  J    Randomize input order of sequences?  No. Use input order
  O                        Outgroup root?  No, use as outgroup species  1
  M        Analyze multiple data sets?  No
  I        Input sequences interleaved?  Yes
  O    Terminal type (IBM PC, ANSI, none)?  ANSI
  1    Print out the data at start of run  No
  2  Print indications of progress of run  Yes
  3                        Print out tree  Yes
  4        Write out trees onto tree file?  Yes
  5  Reconstruct hypothetical sequences?  No

  Y to accept these or type the letter for one to change
```
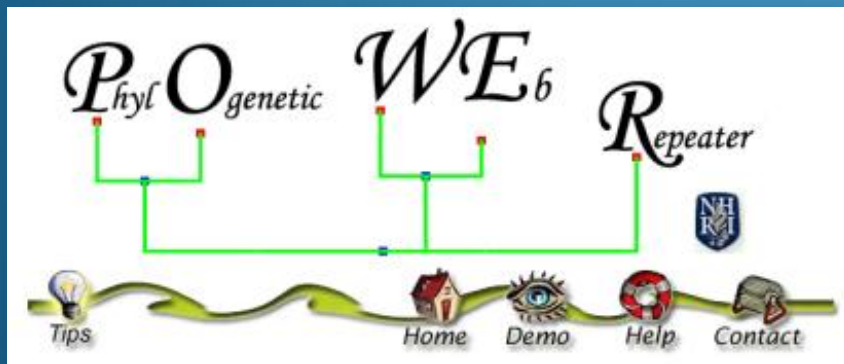
*At this stage they do not have a mouse-windows interface for PHYLIP*

Nucleic Acids Research, 2005

# General Pipeline for Phylogenetic Analysis



Multiple Sequence Alignment

| Methods | Nucleic acid | Protein |
|---|---|---|
| Character state methods | • Maximum parsimony (heuristic search) method<br>• Maximum parsimony (branch and bound search) method<br>• Compatibility method | • Maximum parsimony (heuristic search) method |
| Distance Methods | • Distance matrix computation<br>• Neighbor-joining and UPGMA method<br>• Fitch-Margoliash and least squares method<br>• Fitch-Margoliash and least squares method with molecular clock | • Distance matrix computation<br>• Neighbor-joining and UPGMA method<br>• Fitch-Margoliash and least squares method<br>• Fitch-Margoliash and least squares method with molecular clock |
| Maximum likelihood mothodes | • Maximum likelihood method<br>• Maximum likelihood method with molecular clock | |

Selection of inference Methods

Substitution Model
Tree Construction

Bootstrapping

Evaluate phylogenetic tree

# *Flowchart of Analysis*



(Mount, *Bioinformatics*)

# *Phylogenetic Analysis Tool*

# POWER: PhylOgenetic WEb Repeater

➢ Provide a seamless way to conduct the complex phylogenetic analysis for Biologists

➢ An integrated and user-optimized framework for biomolecular phylogenetic analysis

➢ POWER uses an open-source LAMP (Linux, Apache, MySQL, PHP) structure and infers genetic distances and phylogenetic relationships using well-established algorithms (ClustalW and PHYLIP)

➢ Through a user-friendly web interface, users can sketch a tree effortlessly in multiple steps

➢ Furthermore, iterative tree construction can be performed by adding sequences to, or removing them from, a previously submitted job

# *Integration of Phylip Packages into Automatic Flow*

# *Inside of POWER*

# *POWER: PhylOgenetic WEb Repeater*

## *http://power.nhri.org.tw*

The PhylOgenetic Web Repeater (POWER) allows users performing phylogenetic analysis with molecular data by most programs of PHYLIP package repeatedly. POWER provide two pipelines to process the analysis. One of them includes multiple sequence alignment (MSA) at the begining of the pipeline whereas the other begin phylogenetic analysis with aligned sequence.

Please start your analysis by selecting the pipeline and the data type:

| Pipeline | ⦿ MSA + Phylogenetic Analysis(Input the FASTA format) |
| | ○ Phylogenetic Analysis Only(Input the PHYLIP format) |
| Sequence Type | ⦿ DNA |
| | ○ Protein |

# PhylOgenetic Web Repeater (POWER)



*Data Input*

*MSA parameter selection*

*Phylogeny inference*

# PhylOgenetic Web Repeater (POWER)

**Options of bootstrapping**



**Selection of substitution model**



**Selected method for phylogeny inference**

# PhylOgenetic Web Repeater (POWER)

## Result and Logs

*Online or as bookmark*



*Or E-mail notification*





*Re-perform the process by items added or deleted*

# *PhylOgenetic Web Repeater (POWER)*



Add/ delete sequences to invoke new job

# *Publication in POWER*



The Chloroplast Genome of *Phalaenopsis aphrodite* (Orchidaceae): Comparative Analysis of Evolutionary Rate with that of Grasses and Its Phylogenetic Implications

*Mol. Biol. Evol. 23(2):279–291. 2006*

Ching-Chun Chang,*[1] Hsien-Chia Lin,*[1] I-Pin Lin,† Teh-Yuan Chow,‡[2] Hong-Hwa Chen,* Wen-Huei Chen,§ Chia-Hsiung Cheng,‡ Chung-Yen Lin,|| Shu-Mei Liu,‡ Chien-Chang Chang,¶ and Shu-Miaw Chaw¶

*Institute of Biotechnology, National Cheng Kung University, Tainan, Taiwan; †Department of Superintendent, Tainan Municipal Hospital, Tainan, Taiwan; ‡Institute of Plant and Microbial Biology, Academia Sinica, Taipei, Taiwan; §Department of Life Sciences, National University of Kaohsiung, Kaohsiung, Taiwan; ||Institute of Information Science, Academia Sinica, Taipei, Taiwan; and ¶Research Center for Biodiversity, Academia Sinica, Taipei, Taiwan

# *Service Usage of POWER from 2005 July.*



Accumulative Visit by Country

Accumulative Sequences by Country

# Service Usage of POWER from 2005 July.



Sequence (Accumulative)

More than 228,000 sequences



Visit (Accumulative)

Near 11,800 Visits

# *Automatic Online Demonstration*



http://power.nhri.org.tw/, in the Demo page

# *Conduct Distance Method in POWER*

*Vibrionaceae* dominates the microflora antagonistic towards *Listonella anguillarum* in the intestine of cultured Atlantic cod (*Gadus morhua* L.) larvae

Anders Jón Fjellheim [a], Karina Jane Playfoot [a], Jorunn Skjermo [b,*], Olav Vadstein [c]

[a] Brattøra Research Center, Department of Biology, Norwegian University of Science and Technology (NTNU), 7491 Trondheim, Norway
[b] SINTEF Fisheries and Aquaculture, Department of Marine Resources Technology, 7465 Trondheim, Norway
[c] Department of Biotechnology, Norwegian University of Science and Technology (NTNU), 7491 Trondheim, Norway

The DNA sequences were aligned to known sequences in the GenBank database using BLAST (Altschul et al., 1990). Phylogenetic relationships were inferred using the neighbour joining method (NJ), based on the Kimura two-parameter model (K2P), in **the PhylOgenetic WEb Repeater (POWER)** (Lin et al., 2005).

# NJ with 1000 Replicates in POWER



Phylogenetic tree of urea transporters across taxa. The protein sequences were aligned using ClustalW software, followed by neighbour-joining (NJ) matrix for tree reconstruction and evaluated by means of a bootstrap of 1000 replicates at http://power.nhri.org.tw

# *Perform ML in POWER*

Characterization of a 1-aminocyclopropane-1-carboxylate synthase gene from loblolly pine (*Pinus taeda* L.)

J.R. Barnes [a,1], W.W. Lorenz [b], J.F.D. Dean [b,*]

[a] Department of Biochemistry and Molecular Biology, University of Georgia, Athens, GA 30602, USA
[b] Warnell School of Forestry and Natural Resources, University of Georgia, Athens, GA 30602, USA

Phylogenetic tree depicting the relatedness of ACC synthase and aminotransferase protein sequences in GenBank. The phylogenetic tree was generated using the **POWER server (http://power.nhri.org.tw/)** with default parameters for the maximum likelihood method (ML) and molecular clock, but without bootstrapping or random input of sequences.

# *Execute MP in POWER*

## Molecular characterization and quantification of the gonadotropin receptors FSH-R and LH-R from Atlantic cod (*Gadus morhua*)

C. Mittelholzer [a,*], E. Andersson [b], G.L. Taranger [b], D. Consten [a,1], T. Hirai [c], B. Senthilkumaran [d], Y. Nagahama [e], B. Norberg [a]

[a] Institute of Marine Research Austevoll, N-5392 Storebø, Norway
[b] Institute of Marine Research, N-5817 Bergen, Norway
[c] Department of Biosciences, Teikyo University of Science and Technology, Uenohara, Yamanashi 409-0193, Japan
[d] Department of Animal Sciences, School of Life Sciences, University of Hyderabad, Hyderabad 500 046, India
[e] Laboratory of Reproductive Biology, National Institute for Basic Biology, 444-8585 Okazaki, Japan

Phylogenetic comparison of fish full-length FSH-R and LH-R amino acid sequences analysed by **POWER** using **maximum parsimony (MP) and default settings**. A rooted consensus phylogenetic tree generated by means of the Neighbor-Joining algorithm, using the LGR sequence of the fruit fly (Drosophila melanogaster) and sea lamprey (Petromyzon marinus) as outgroups was drawn with njplot. Bootstrap values from 1000 replicates are indicated for each tree node.

# *POWER Listed in*

- PHYLIP Programs maintained by Joe Felsenstein
  - Recent listings:
    - POWER server (26 August 2007) to align sequences and infer phylogenies, http://evolution.genetics.washington.edu/phylip/software.serv.html
- BioToolKit by CSHL press (BioSupplynet.com)
    - ALL CATEGORIES / GENOMICS RESOURCES / EVOLUTIONARY AND COMPARATIVE BIOLOGY (80)
- Bioinformatics Links Directory
    - DNA : Phylogeny Reconstruction
- ONLINE ANALYSIS TOOLS (http://molbiol-tools.ca/)
- ExPASy (Phylogenetics and taxonomy databases & resources)



Phylogenetics and taxonomy databases & resources
- COG - Phylogenetic classification of proteins encoded in complete genomes
- EGO - Eukaryotic Gene Orthologs
- InParanoid - Eukaryotic ortholog groups
- Metazome - Phylogenomic analysis of metazoan gene families
- OMA - Orthologs Matrix Project (OMA)
- TreeBASE - Relational db of phylogenetic information
- TreeFam - Tree families database of phylogenetic trees of animal genes
- The PhylOgenetic Web Repeater (POWER) - perform phylogenetic analysis
- NEWT - UniProt Taxonomy Browser
- CluSTr - Automatic classification of UniProtKB proteins into groups of related proteins
- ProtoNet - Classification of the proteins into hierarchical clusters

# Distance Method, MP and ML

- Which method should we choose?
- The main disadvantage of distance-matrix methods is their inability to efficiently use information about local high-variation regions that appear across multiple subtrees.
- ML is broadly similar to the maximum-parsimony (MP) method, but maximum likelihood allows additional statistical flexibility by permitting varying rates of evolution across both lineages and sites.
- ML, a better choice?

# *Maximum Likelihood*

- Conditional probability of the data (Aligned sequences) given a hypothesis (a model of substitution with a set of parameter ө, and the tree τ, including topology and branch lengths)

$$L(\tau, ө )=\text{Prob}(\text{Data}| \tau, ө )$$

Or

Prob(Aligned Sequences| tree, model of evolution)

# *Relationships among some standard models of nucleotide evolution*

# Illustration of DNA Substitution Model

$$Q = \begin{pmatrix} -(x_1 + x_2 + x_3) & x_1 & x_2 & x_3 \\ \frac{\pi_1 x_1}{\pi_2} & -\left(\frac{\pi_1 x_1}{\pi_2} + x_4 + x_5\right) & x_4 & x_5 \\ \frac{\pi_1 x_2}{\pi_3} & \frac{\pi_2 x_4}{\pi_3} & -\left(\frac{\pi_1 x_2}{\pi_3} + \frac{\pi_2 x_4}{\pi_3} + x_6\right) & x_6 \\ \frac{\pi_1 x_3}{\pi_4} & \frac{\pi_2 x_5}{\pi_4} & \frac{\pi_3 x_6}{\pi_4} & -\left(\frac{\pi_1 x_3}{\pi_4} + \frac{\pi_2 x_5}{\pi_4} + \frac{\pi_3 x_6}{\pi_4}\right) \end{pmatrix}$$



Base frequencies   ○  ○  ○  ○
                       C  G  T
                  A  ☐  ☐  ☐
Substitution rates   C  ☐  ☐
                  G  ☐

GTR (for four characters, as is often the case in phylogenetics) requires 6 substitution rate parameters (x1~x6), as well as 4 equilibrium base frequency parameters.



GTR

# *Illustration of Models for DNA*

# *Models of Amino Acid Replacement*

*Phylogenetic Reconstruction by Automatic Likelihood Model Selector (PALM) :*
*A Framework for Phylogenetic Analysis with the Best Substitution Model*

陳淑華
sophia@iis.sinica.edu.tw

PLoS ONE, 2009

# *Background for PALM*

- Likelihood methods in phylogenetics relaxes the parameters for varying rates of evolution across both lineages and sites, which is robust in dealing with various extend of input sequence similarity.

- Model fitting has been suggested for many years, but many researchers select models arbitrarily. They often feel confusing either in making choice among models, or in dealing with the conflict on the results concluded by different models.

- The computing of likelihood method is intensive. Thus the ML-based model selecting procedure is hard to implement.

- Here we present the way to identify the best-fit model based on liklihood measurement. Consequently, model fitting is possible to be a routine practice integrated in a phylogenetic analysis.

# *Motivation I*

➤ Provide a seamless way to conduct the complicated phylogenetic analysis for biologists and biomedical researchers.

➤ An integrated and user-friendly framework for conducting molecular phylogenetic analysis

➤ PALM is constructed on an open-source LAPP (Linux, Apache, PostgreSql, PHP) structure

➤ PALM infers genetic distances and phylogenetic relationships using well-established algorithms (ClustalW , PhyML, ProtTest, Modeltest) in an automatic pipeline.

# *Motivation II*

➢ Fitness of model can be measured and selected by following criteria: likelihood ratio tests (hLRTs), Akaike information criterion (AIC), and Bayesian information criterion (BIC)

➢ PALM helps user to construct the phylogenetic relationship by ML-based method with bootstrap using the best-fit substitution model.

➢ Through the friendly web interface, users can sketch a phylogenetic tree effortlessly

➢ Furthermore, iteration on phylogenetic reconstruction is possible by adding sequences to, or removing them from a previously result.

# *Component Programs of PALM*

➤ PhyML 3.0

➤ ModelTest 3.7

➤ ProtTest 2.0

➤ ClustalW 2.0.3

➤ ReadSeq

# *Models Used in PALM*

- For DNA (56 models)

  - JC69, K80, F81, HKY, TrN, TrNef, K3P, K3Puf, TIM, TIMef, TVM, TVMef, SYM, GTR
  - Options of +I, +G

- For Protein (112 models), **Time consuming**

  - LG, DCMut, JTT, MtREV, MtMam, MtArt, Dayhoff, WAG, RtREV, CpREV, Blosum62, VT, HIVb, HIVw
  - Options of +I, +G, +F

# *Flowchart of PALM*



> Seq1
AAAATTTC...
>Seq2
AATTCGGAC..
.......

Sequences in Fasta

**Input**

FASTA Sequence

alignment

Aligned Sequence

**Input**

Aligned sequences

*DNA*

*Protein*

*PhyML/ Modeltest/ **PALMmonitor*** (56 Models)

*PhyML/ Prottest/ **PALMmonitor*** (112 Models)

Add/ delete nodes & re-submit input

Controlled by **PALM Daemon**

*Best Model selected by AIC/ AICc/ BIC/ LnL*

bootstrapping

*Bootstrap Result & Tree Reconstruction by PALMtree*

*PhyML-MPICH2/ PALM Parallel Job controller*

# Distribution Computing by PalmMonitor for the Likelihood Estimation of Models

56 for DNA/ 112 for Protein

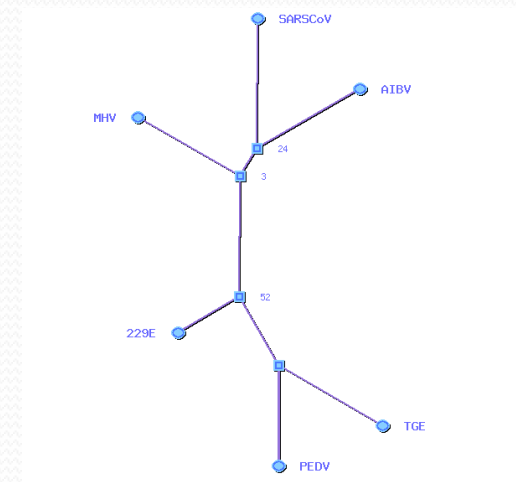7 parallel processes

n   n   n   n   n   n   n

# *Input and Output of PALM*

- Input format (Protein and DNA)
  - FASTA format
  - Phylip format: Aligned Sequences
  - User tree (if a valid tree is submitted)
- Output
  - Tree topology
  - Tree file in Newick format
  - Aligned sequence in phylip format
  - The best model selected by PALM
  - Likelihoods of all available models

# *Result of PALM*



**A** — PALM logo and PALM Result header

| Job ID | 20080821060606361 | Number of Substitution Rate Category | 4 |
|---|---|---|---|
| Job Note | test for speed in protein | Model Selection Criterion | LnL |
| Sequence Type | Protein | Optimization of Tree Topology | Yes |
| Number of Bootstrap | 1000 | Optimization of Branch Length | Yes |
| Starting Tree | BIONJ | | |

**C**

| Best Model Selected | JTT+I+G+F |
|---|---|
| Model Selection Criterion | LnL |
| AIC | 2336.50 |
| -lnL | 1134.25 |

**D**

| Model | deltaAIC | AIC | -lnL* | AICw |
|---|---|---|---|---|
| JTT+I+G+F | 2.00 | 2336.50 | -1134.25 | 0.12 |
| JTT+G+F | 0.00 | 2334.50 | -1134.25 | 0.33 |
| WAG+I+G+F | 2.78 | 2337.28 | -1134.64 | 0.08 |
| WAG+G+F | 0.78 | 2335.28 | -1134.64 | 0.23 |
| WAG+I+F | 5.17 | 2339.67 | -1136.83 | 0.03 |
| WAG+F | 3.50 | 2338.00 | -1137.00 | 0.06 |

**B** — Phylogenetic tree diagram

**E**

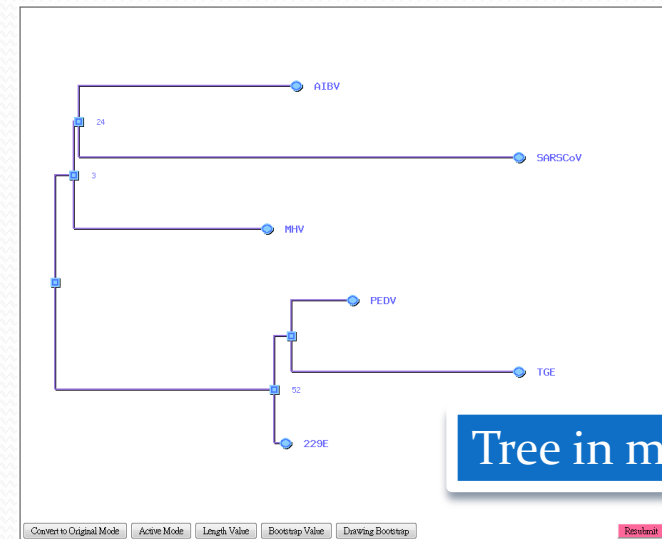| Original File | 20080821060606361 |
|---|---|
| Phylip File | 20080821060606361.phy |
| Phylogenetic Tree (Newick) | tree20080821060606361.txt |
| ProtTest Information | ProtTest_20080821060606361.txt |
| Bootstrap Tree | 20080821060606361_phyml_boot_trees.txt |
| Bootstrap Statistic data | 20080821060606361_phyml_boot_stats.txt |

The job is computed approximately in 47 minute(s).

# *PALMtree*



Unroot Tree

Original Tree

Tree in midpoint

# *Demo Flash of PALM*



*http://palm.iis.sinica.edu.tw/demo.html*

# *Some Suggestions*

- Please be patient and make a reasonable choice of the input sequence set

- Only **well aligned** sequences lead to meaningful phylogenetic result.

- RNA editing may introduce bias during analysis. Avoid those regions that may have such conditions.

# *Bootstrap (BS) Analysis*

- Bootstrap analysis is the most popular method for statistical evaluation of phylogenies.

- In general:

  - **BS >95%: Often close to 100% confidence in that branch**

  - **BS>75%: Often close to 95% confidence in that branch**

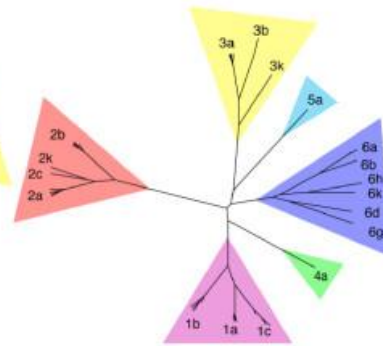  - BS<75% : Maybe a correct clade, while the original bias cannot be corrected by the re-sampling process.

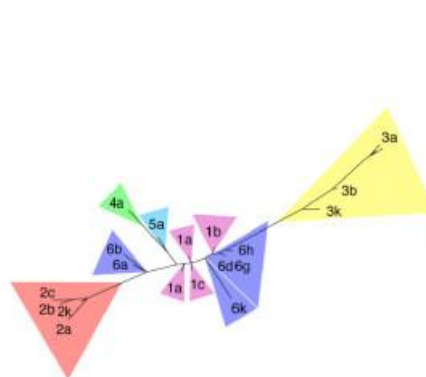# *Input Sequences Make the Tree Different*
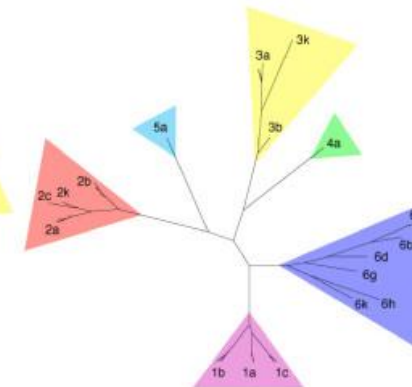
HIV



(a) Complete Genome
(b) Polyprotein
(c) 5' UTR
(d) Okamoto region of NS5B

# *Future Plans for PALM*

- Gateway to integrate users-defined substitution models
- Stand-a-lone version of *PALM*
- Improve and optimize the performance of whole pipeline by applying <span style="color:red">parallel computing/cloud computing</span>
- Implement of advanced, sophisticated phylogenetic inference methods such as MrBayes.

# *Acknowledgement*

National Health Research Institutes
國家衛生研究院

Chia-Ling Chen
Chieh-Hwa Lin
Li-Wei Lai
Shu-Juan Hsu
Ming-Hsin Tasi
Chao A. Hsiung

中央研究院
資訊科學研究所
Institute of Information Science
Academia Sinica

Daniel, Sheng-Yao, Su
Pan-Han Kuo
Tengi Huang
Chen-Zen Lo
Linda, Yi-Shuian Lu

# Bioinformatics Core for Genomic Medicine and Biotechnology Development



*http://www.tbi.org.tw*