# Deciphering the Biological Problems in the Approach of Systems Biology

林 仲 彥, *Chung-Yen Lin Ph.D.*

*Assistant Research Fellow*

*cylin@iis.sinica.edu.tw*

*Laboratory of Systems Biology and Network Biology*
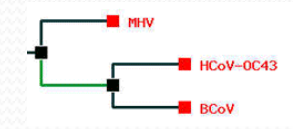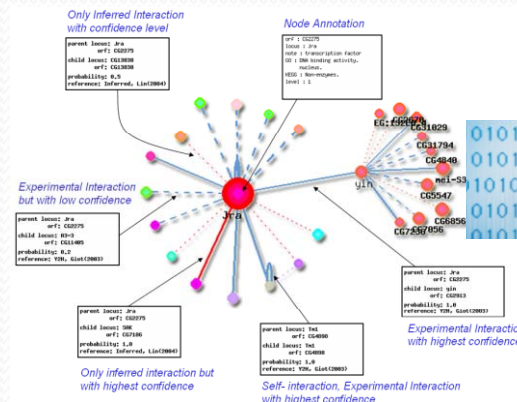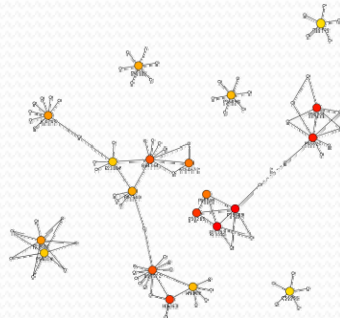
*Institute of Information Science, Academia Sinica*

*Aug 12, 2008*

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

NHRI 國家衛生研究院
National Health Research Institutes

# *Outline*

- **Solving the Biological problems in Computational methods and statistics**
  - Genomics studies for high throughput research
  - Phylogenetics analysis
  - Protein interactome
  - Network comparison and topological analysis
  - Ongoing projects

# *Platform Based on LAMP/ LAPP*

**L**inux

    Operation System

**A**pache (with OpenSSL)

    Webserver

**M**ySQL/**P**ostgrSQL
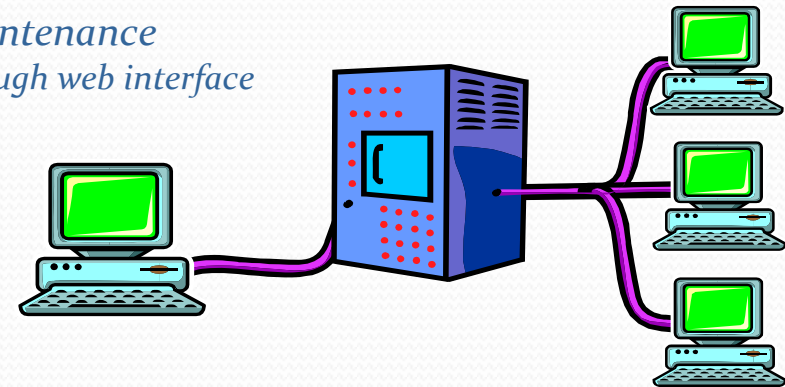
    Relational Database

**P**HP

    Server-side HTML embedded
      scripting language with GD
      library

*Query by Web interface*
*(Win9x/Me/2000/Mac/Unix*
*/Linux/Solaris)*

*Maintenance*
*through web interface*

*Web database*
*Linux*
*Apache webServer,*
*MySQL*
*PHP control language*

# Design of Primers and Probes
## for High and low Throughput Research

- ✓ Primer Design Assistant (PDA)
- ✓ Unique Probe Selector (UPS)

# *Motivation for PDA*

- **Integration of experiences from wet-lab and computational technology** to perform primer design in large scale PCR under similar Tm value.

- Primer Design Assistant (PDA) is a **web interface primer design service combined with thermodynamic theory** to evaluate the fitness of primers.

- It runs in a Linux–Apache–MySQL–PHP structure on a PC equipped with dual CPU (Intel Pentium III 1.4 GHz) and 512 Mb of RAM.

- **A succinct user interface** of PDA is accomplished by built-in parameters setting. Advanced options on 5' GC content, 3' GC content, dimer check and hairpin check are available.

- **PDA accepts single sequence query or multiple ones in FASTA format**. It produces optimal and homogenous primer pairs that meet the need in experimental design with **large-scaled PCR** amplifications.

# *Genomics Studies For High Throughput Research : PDA*



- **Primers designed through PDA has been experimentally proved to reach 97% successful rate**
- **PDA can be used to design the primers set for high through put experiments. For example , for 96 /384 format PCR Rx.**
- **http://dbb.nhri.org.tw/primer/**
- **Published on NAR 2003**

# *Criterion for PDA Setting*

- Default Settings

- Advanced options

| Repeats | Any four continual nucleotides (AAAA, TTTT, CCCC, or GGGG) will be avoided for both forward and reversed primers. Continuous dinucleotide repeats, such as 'ATATAT', are also avoided. |
|---------|--------|
| C/G clamp | G or C on the end of 3' terminal |
| GC % | 25% ~ 75% |
| Tm | Tm of forward and reversed primers restricted to be higher than 50°C |
| ΔTm | restricted to be smaller than 5°C |

| Dimer check: | This option turns on can avoid primer dimer formation. |
|--------------|--------|
| Hairpin check: | This option turns on can avoid internal self-complementarity. |
| 5' GC content check: | Check the GC% of 5' to add the ability to recognize the template and enhance the priming specificity. |
| 3' GC content check: | Check the GC% of 3' to avoid mismatch to avoid mismatch. |
| Covered region: | By entering the start position and stop position, you can get the PCR product containing the segment you need. |

$$T_{\mathrm{m}}(°\mathrm{C}) = 59.9 + (0.41 \times \text{GC content}) - \left( \frac{675}{\text{primer length}} \right)$$

# Calculation of the Stability of DNA Duplexes

**A** Primer-to-primer annealing

primer 1  5' [atcgt]  3'
primer 2  3' [tagcc]  5'

sliding direction

**B** Hairpin structure

5' atcgt
3' tagcc

sliding direction

**C** Primer-to-template annealing

primer
template

**D** Nearest-neighbor parameters for all possible NN dimer duplexes. Modified from SantaLucia 1998 (10).

| Sequence | Free Energy Parameter ($\Delta G^{\circ}_{56}$) |
|---|---|
| A a | -0.73 |
| A t | -0.61 |
| A c | -1.16 |
| A g | -0.82 |
| T a | -0.32 |
| T t | -0.73 |
| T c | -1.03 |
| T g | -1.16 |
| C a | -1.16 |
| C t | -0.82 |
| C c | -1.57 |
| C g | -1.81 |
| G a | -1.03 |
| G t | -1.16 |
| G c | -1.92 |
| G g | -1.57 |
| A | 0.98 |
| T | 0.98 |
| C | 1.00 |
| G | 1.00 |

5'- ATCGT -3'
****
3'- TAGCC- 5'

(-0.61)+(-1.03)+(-1.81)+(0.98)= -2.47

When sequence 1 [5'- atcgt -3'] aligns to sequence 2 [3'- tagcc- 5'], the first base of the first sequence (a) matches to (t) in sequence 2, and follows with three more Watson–Crick pairs. The fifth base mismatches. The NN propagation energy of the continuing base pairs (at), (tc), (cg) and the mismatched base (t) in primer 1 are summed up: (-0.61)+(-1.03)+(-1.81)+(0.98)=-2.47.

# *Ranking Mechanism*

The primer pairs passing through the limitations listed above are sorted by ranking score ($R$):

$$R = 100 - \Delta(T_\mathrm{m}) + \Delta G^{\circ}_{\mathrm{forward}}(3' - 5') + \Delta G^{\circ}_{\mathrm{reverse}}(3' - 5') + \mathrm{hairpin\ score} + \mathrm{dimer\ score}$$

To avoid the mis-priming amplification, the 5' end of the primer is expected to anneal to target templates more stable than the 3' end

# *Currently available service (conti)*

➤ **Primer Design Assistant (PDA)**

- Customized PCR conditions

Dimer check  →

Hairpin check  →

5'GC content check  →

3'GC content check  →

Covered region  →

| Input format: | ⦿ fasta  ○ text |
|---|---|
| Sequence(s) input or file upload | |
| | 瀏覽... |
| Primer length: | 19 ▾ |
| PCR product size: | 150 ▾ |
| **Advanced Options** | |
| Dimer check: | ⦿ No  ○ Yes |
| Hairpin check: | ⦿ No  ○ Yes |
| 5' GC content check: | ⦿ No  ○ Yes |
| 3' GC content check: | ⦿ No  ○ Yes |
| Covered region: | Start from [    ] -- End on [    ] |

search       reset

# *Currently available service (conti)*



Batch primer design for unified experimental conditions

| Input format: | ⦿ fasta ○ text |
|---|---|
| Sequence(s) input or file upload | |
| | 瀏覽... |
| Primer length: | 19 |
| PCR product size: | 150 |

# Currently available service (conti)

| | | | criteria | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| format | primer_length | primer_window_size | Repeats avoid ( AAAA, TTTT... ) | C/G clamp | GC % | Tm | △ Tm | dimer check | hairpin check | 5'GC content check | 3'GC content check | covered region | sequence |
| text | 19 | 150 | yes | yes | 25% ~ 75% | ≧ 50℃ | ≦ 5℃ | no | no | no | no | ~ | atggcgtctccttctagaaa... |

| full text | primer | GC% | Tm | offset | rank | PCR product |
|---|---|---|---|---|---|---|
| forward primer | cccgctgttcaccctgttc | 63.16 | 50.27 | 3530 | 1 | cccgctgttcaccctgttct... |
| reverse primer | ccggaccctgaccaaatcc | 63.16 | 50.27 | 3679 | | |
| | | | | | | |
| forward primer | ggcaggccgagcaattcag | 63.16 | 50.27 | 728 | 2 | ggcaggccgagcaattcagt |

**PDA Report page**

**Convenient Excel download format**

| | E | F | G | H | I | J | K | L | M | N |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | C/G clamp | GC content | Tm | △Tm | dimer chec | hairpin che | 5' GC chec | 3' GC chec | covered region | |
| 2 | yes | 25% ~ 75% | ≧50℃ | ≦5℃ | no | no | no | no | ~ | |
| 3 | | | | | | | | | | |
| 4 | | | | | | | | | | |
| 5 | | | | | | | | | | |
| 6 | aggagccgctccgatacagctacaaccccgaccagttccacaacatggacctcaggggcggcccccacgatggcgtcaccattccccgctccaccagcgac |
| 7 | | | | | | | | | | |
| 8 | offset | rank | PCR product | | | | | | | |
| 9 | 3530 | 1 | cccgctgttcaccctgttcttcacctcagaactccccaggtttacagagagccagtgcaagagcccttccccctaccgaaga |
| 10 | 3679 | | | | | | | | | |
| 11 | | | | | | | | | | |

A1 = format

| forward primer | ccagccacacgcgc... |
|---|---|
| reverse primer | tcgtggactccgggt... |
| | |

# *Primer set for Nested PCR in PDA*

Input target sequence



Get the primer set and their location

# Primer set for Nested PCR in PDA

Modify the size of Product and fill the location of 1[st] PCR product

Get the primer sets for nested PCR with 500 bps product



Input format: ◉ fasta ○ text

Sequence(s) input or file upload
>AB048365.1:21..4778
atggcgtctccttctagaaactcccagagccgacgccggtgcaagg
agccgctccgatacagctacaaccccgaccagttccacaacatgg
acctcagggcgcaacccccacgatggcgtcaccattcccagctccac

Primer length: 19

PCR product size: 500

Advanced Options

Dimer check: ◉ No ○ Yes

Hairpin check: ◉ No ○ Yes

5' GC content check: ◉ No ○ Yes

3' GC content check: ◉ No ○ Yes

Covered region: Start from 3530 -- End on 3679

| format | primer_length | primer_window_size | Repeats avoid ( AAAA, TTTT... ) | C/G clamp | GC % | Tm | △ Tm |
|--------|---------------|--------------------|--------------------------------|-----------|------|----|------|
| fasta | 19 | 500 | yes | yes | 25% ~ 75% | ≥50℃ | ≤5℃ |

| partial text | primer | GC% | Tm | offset | rank | |
|--------------|--------|-----|-----|--------|------|--|
| forward primer | cctgcaggctgccttccac | 68.42 | 52.43 | 3492 | 1 | cctgcaggctgcc |
| reverse primer | gagccagacccaggatgcg | 68.42 | 52.43 | 3991 | | |
| | | | | | | |
| forward primer | tgcaggctgccttccaccc | 68.42 | 52.43 | 3494 | 2 | tgcaggctgcctt |
| reverse primer | cagagccagacccaggatg | 63.16 | 50.27 | 3993 | | |

# *Use PDA to Develop PCR kits for SARS Detection*

➤ 由於2003年當時通用的SARS-CoV檢驗方法靈敏度有限，使得病毒量較低或是因採樣方法不佳的檢體無法被檢測到，在防疫的前提下，本核心便與國家衛生研究院基因醫學研究組協同疾病管制局(CDC, Taiwan)發展出高靈敏度之檢測方法。

➤ 檢測方法中所需要的核酸引子都透過PDA來進行設計，避免引子本身dimer 及 hairpins 的形成，加速了檢測方法的建立。

➤ 此方法為結合1st run RT-PCR + 2nd run Q-PCR，可於1.5小時內檢測出結果，經實驗證明縱使病毒量低於10隻，也可以透過這一套方法檢測出來。

# *Result of Real-time PCR for SARS-CoV Detection*

# Primer Design Assistant (PDA)



From July 2003 to Aug 2008

# *Primer Design Assistant (PDA)*



*From July 2003 to Aug 2008*
*Over 100,000 Visits and 600,000 submitted Seqs*

Nucleic Acids Res. 2003, 31: 3751-3754

# *Usage of PDA Worldwide*



Visits

Submitted Seqs

Sequences submitted to PDA from overseas, accumulative,
Jul. 2003 to Apr. 2008; with 334,426 sequences submitted from Taiwan

# BMC Bioinformatics

Software

## ProbeMaker: an extensible framework for design of sets of oligonucleotide probes

Johan Stenberg*, Mats Nilsson and Ulf Landegren

Address: Department of Genetics and Pathology, Rudbeck Laboratory, Uppsala University, Se-751 85, Uppsala, Sweden

Email: Johan Stenberg* - johan.stenberg@genpat.uu.se; Mats Nilsson - mats.nilsson@genpat.uu.se; Ulf Landegren - ulf.landegren@genpat.uu.se

* Corresponding author

This article is available from: http://www.biomedcentral.com/1471-2105/6/229

*BMC Bioinformatics* 2005, **6**:229

11. Chen SH, Lin CY, Cho CS, Lo CZ, Hsiung CA: **Primer Design Assistant (PDA): A web-based primer design tool.** *Nucleic Acids Res* 2003, **31**:3751-3754.
12. Vallone PM, Butler JM: **AutoDimer: a screening tool for primer-dimer and hairpin structures.** *Biotechniques* 2004, **37**:226-231.
13. Kaderali L, Schliep A: **Selecting signature oligonucleotides to identify organisms using DNA arrays.** *Bioinformatics* 2002, **18**:1340-1349.
14. Rouillard JM, Zuker M, Gulari E: **OligoArray 2.0: design of oligonucleotide probes for DNA microarrays using a thermodynamic approach.** *Nucleic Acids Res* 2003, **31**:3057-3062.
15. Kaderali L, Deshpande A, Nolan JP, White PS: **Primer-design for multiplexed genotyping.** *Nucleic Acids Res* 2003, **31**:1796-1802.
16. Shoemaker DD, Lashkari DA, Morris D, Mittmann M, Davis RW: **Quantitative phenotypic analysis of yeast deletion mutants using a highly parallel molecular bar-coding strategy.** *Nat Genet* 1996, **14**:450-456.
17. Stenberg J: **The MolTools Java library.** [http://sourceforge.net/

Chen SH, Lin CY, Cho CS, Lo CZ, Hsiung CA:  Primer Design Assistant (PDA): A web-based primer design tool.  Nucleic Acids Res 2003, 31:3751-3754.

PDA
Primer Design Assistant

Help    Contact

Input format:        fasta    text
Sequence(s) input
or file upload

Primer length:       19
PCR product size:    150

Chen, S. H., C. Y. Lin, C. S. Cho, C. Z. Lo, and C. A. Hsiung. 2003. Primer
Design Assistant (PDA): a web-based primer design tool. Nucleic Acids Res.
31:3751–3754

# Genomics Studies For High Throughput Research : UPS

➢ **Unique Probe Selector (UPS)**

➢ *Probe design for hybridization in low and high throughput experiments*

➢ **http://array.iis.sinica.edu.tw/ups/**

➢ **BMC Bioinformatics 2008**

# Brief on Unique Probe Selector (UPS)

- Although most of these tools aimed on designing probes for microarray, only few of them take the genetic background noise in hybridization reaction into account. **A web tool** for customized probe design regarding to the discriminating power of probes is sparse.

- Here we present a web tool, the Unique Probe Selector (UPS), for selecting unique oligo-nucleotide probe. The algorithms applied here include **thermodynamic theory**, **GC content**, **GC clamps**, **secondary structure of probes** and **some other empirical preferences of wet-lab researchers**. **Low-complexity** regions are filtered out to maintain probe specificity.

- The UPS evaluates probe-to-target hybridization under a user-defined condition *in silico* **to ensure high-performance hybridization** and **minimizes the possibility of non-specific reactions**.

- UPS has been applied to design human arrays for gene expression studies and to develop several small arrays of gene families that were inferred as molecular signatures of cancer typing/staging or pathogen signatures.

# *Infrastructure of UPS*

- Under the consideration of efficiency and performance, we adopted the LAMP structure (Mandrake 2007, the operating system , Apache (webserver), PostgreSQL (database) and PHP), to provide the web access, files upload/ download, mail notification and data storage.

- All the calculations related with unique probe design was performed on a window-based machine in Delphi code.

# Basic Criteria for Probe Selection

①    Probe length: from 30 ~ 120 bps

②    Melting Temperature

    ✓  The probe annealing temperature (Ta) is determined by melting temperature (Tm).  Probe Tm depends on several physiochemical factors and is calculated in the following equation based on Nearest Neighbor model

$$\Delta Tm = \Delta H / (10.8 + \Delta S + R \times \ln (C / 4)) - 273.15 + 16.6(\log_{10}[Salt])$$

③    Sequence complexity

    ✓  We exclude any five or more continual nucleotides (AAAAA, TTTTT, CCCCC, or GGGGG).  Continuous di-nucleotide/ tri-nucleotide repeats, such as 'ATATAT" and 'ATGATGATG', are also avoided.

# *Basic Criteria for Probe Selection*

④   Computation of secondary structure formation

- ✓ We use a perl program UNAFold.pl integrated into UPS to calculate ΔG

⑤   Continuous stretch and identity between probe and no-target template

- ✓ Here we used Li *et al's* (NAR 2003) experimentally established criteria to exclude unsuitable oligonucleotides: **identity of ≧ 85%, continuous stretch of ≧17 and free energy <-35 kcal/mol** (it will depend on the length of probe) between probe and non-target templates.

# Demonstration of UPS 2.0

# Demonstration of UPS 2.0

| | |
|---|---|
| **Probe Uniqueness*** | ⦿ Unique Probe within group |
| | ○ Unique Probe based on the specific organism<br>Aedes_aegypti (yellow fever mosquito) ▾ |
| | ○ Unique Probe based on user's defined organism<br>[ ] 瀏覽… |
| Sequence (s) Paste or File upload* | >ABL<br>AAGGTAGCTGATTTTGGCCTGAGCAGGTTGATGA<br>CAGGGGACACCTACACAGCCCATGCTGGAGCCA<br>AGTTCCCCA<br>[ ] 瀏覽… **DEMO** |
| Probe Length | 70 ▾ |
| Probe # for each sequence | 1 ▾ (maximum 3) |
| Job note (optional) | PTP family |
| E-mail* | yamatolin@gmail.com |
| **Advanced Options** | |
| [Salt] | salt_conc 0.58 (0~1M) |
| Degenerate probe allowed | ○ yes ⦿ no |

submit reset

# *Jobs Accepted by UPS*



**UPS** Unique Probe Selector

Home   Demo   Help   Contact

Dear Sir,

We accepted your submission. The job will be done in a few minutes to hours. After the job being finished, you will receive a notice email, or you can check the result from the link below.

http://array.iis.sinica.edu.tw/ups/result.php?ID=20070827232034

[ Add to my favorite ]

Thanks for using UPS. Any comment will be appreciated.

Your faithfully.

UPS Administrator.

# *Notification by Email*

**Message from UPS , time stamp : 2007/08/27 - 23:21:22**

☆ **UPS administrator** 寄給 我

Dear Sir or Madam,

The job 'PTP family' you sent has finished!

You can check the result from the link below.

Thank you for using UPS.

Your faithfully.

UPS Administrator.

--------------

Job ID : 20070827232034

http://array.iis.sinica.edu.tw/ups/result.php?ID=20070827232034

May the UPS with you.

# *Output of UPS*

Job Note : PTP family
Type of Probe Uniqueness : Unique Probe within group

Page 1 ⌄

## Output for UPS

**Total : 111**

Advanced Options filter

| Sequence_ID | Rank | CG% | Tm | probe sequence | delta G |
|---|---|---|---|---|---|
| ABL | 1 | 56 | 73 | ctgagcaggttgatgacaggggacacctacacagcccatgctggagccaagttccccatcaaatggactg | 0.183 |
| ARG | 1 | 40 | 68 | gagccaaatttcctattaagtggacagcaccagagagtcttgcctacaataccttctcaattaaatctga | 1.154 |
| EGFR | 1 | 44 | 69 | gcagaaggaggcaaagtgcctatcaagtggatggcattggaatcaattttacacagaatctatacccacc | -1.485 |
| TNK1 | 1 | 76 | 78 | tggtgcggcctctgggcggtgcccggggccgctacgtcatgggcgggcccgccctatcccctacacctg | -3.79 |
| TXK | 1 | 40 | 68 | agccaagttcccaatcaagtggtccctcctgaagttttctttcaataagtacagcagtaaatctgat | 0.802 |
| TYK2 | 1 | 69 | 77 | cctagccaaggccgtgcccgaaggccacgagtactaccgcgtgcgcgaggatggggacagcccccgtgttc | -1.649 |
| TYRO3 | 1 | 54 | 73 | tcggactctcccggaagatctacagtggggactactatcgtcaaggctgtgcctccaaactgcctgtcaa | -0.65 |
| VEGFR1 | 1 | 44 | 69 | gccttgcccgggatatttataagaaccccgattatgtgagaaaaggagatactcgacttcctctgaaatg | 0.096 |
| VEGFR2 | 1 | 44 | 70 | gcccgggatatttataaagatccagattatgtcagaaaaggagatgctcgcctcccttttgaaatggatgg | 0.232 |
| VEGFR3 | 1 | 66 | 75 | gccttgcccgggacatctacaaagaccccgactacgtccgcaagggcagtgcccggctgcccctgaagtg | -1.705 |

## Output for Download

We provide more information for each probe in following files.

1. Best probes in fasta format ⬇

2. All probes in fasta format ⬇

3. All probes in CSV (with Tm, CG%, deltaG, Best_hit, Max_overlap, Identity ) ⬇

4. In silico hybridization check for each probe by BlastN ⬇

# *Advanced Options*

Job No...
Type

Page 1 ▾

**Total : 111**

| Sequence_ID |
| --- |
| ABL |
| ARG |
| EGFR |
| TNK1 |
| TXK |
| TYK2 |
| TYRO3 |
| VEGFR1 |
| VEGFR2 |
| VEGFR3 |

We provid
1. Best prob
2. All probe:
3. All probe:
4. In silico h

| Advanced Options Filter | |
| --- | --- |
| GC% | from 35 ▾ % to 65 ▾ % |
| Tm range | not lower than 45 ▾ °C |

submit  reset

| | |
| --- | --- |
| Job Note | : test in safari |
| Type of Probe Uniqueness | : Unique Probe within group |
| GC% | : from 35 % to 65 % |
| Tm range | : not lower than 50 °C |

Page 1 ▾

## Output for UPS

**Total : 105**

| Sequence_ID | Rank | CG% | Tm | probe sequence | delta G |
| --- | --- | --- | --- | --- | --- |
| ABL | 1 | 56 | 73 | ctgagcaggttgatgacaggggacacctacacagcccatgctggagccaagttccccatcaaatggactg | 0.183 |
| ARG | 1 | 40 | 68 | gagccaaatttcctattaagtggacagcaccagagagtcttgcctacaataccttctcaattaaatctga | 1.154 |
| EGFR | 1 | 44 | 69 | gcagaaggaggcaaagtgcctatcaagtggatggcattggaatcaattttacacagaatctatacccacc | -1.485 |
| TXK | 1 | 40 | 68 | agccaagttcccaatcaagtggtcccctcctgaagtttttcttttcaataagtacagcagtaaatctgat | 0.802 |
| TYRO3 | 1 | 54 | 73 | tcggactctcccggaagatctacagtgggggactactatcgtcaaggctgtgcctccaaactgcctgtcaa | -0.65 |

# *Output for Download*

We provide more information for each probe in following files.

① Best probes in fasta format

② All probes in fasta format

③ All probes in CSV (with Tm, CG%, deltaG, Best_hit, Max_overlap, Identity )

④ In silico hybridization check for each probe by BlastN

# *upQPCR :: PDA+ UPS*

- For specific sequence wanted to identify by Q-PCR, the following steps can be used to get the primer pairs and probe

  ① Submit Sequence to PDA for best primer set with specific region (or select by PDA)

  ② Submit the amplicon to UPS and choose the organism you used to get the best probe for Q-PCR

# Phylogenetic Analysis

- ✓ Phylogenetic Web Repeater (POWER)

- ✓ Phylogenetic reconstruction by Automatic Likelihood Model selector (PALM)

# Coding Characters and Defining Homology



*Classical phylogenetic analysis by Morphology*

*Molecular phylogenetic analysis By Bio-Molecules*

# An Example of Phylogenetic Tree

# *Phylogenetic Tree*

- The tree is composed of nodes connected by branches.



(Andy Vierstraete 1999)

- **node :** a node represents a taxonomic unit.
  - Internal nodes
  - External nodes
- **branch (edge):** defines the relationship between the taxa.
- **branch length :** often represents the number of changes that have occurred in that branch.
- **root :** is the common ancestor of all taxa.
- **distance scale :** scale which represents the number of differences between sequences (e.g. 0.1 means 10 % differences between two sequences)

# *Types of Phylogenetic Tree*

- Branch: define relationship between nodes
  - Branch length: longer branch length, more sequence changes



| Cladogram | Phylogram | Ultrametric tree |
|---|---|---|
| no meaning | genetic change | time |
| Rooted with "SARS" | in substitutions per nucleotide | By assuming a molecular clock |
| **Ex. Pasimony** | **Ex. Neighbor-join, ML** | **Ex. UPGMA** |

# 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.



*( David A. Baum et al., Science 11 November 2005: Vol. 310. no. 5750, pp. 979 – 980)*

# The Ways to Construct the tree

- Distance-matrix methods (Dis)
  - Neighbor-joining
  - Fitch-Margoliash method
  - Using outgroups
- Maximum parsimony (MP)
  - Branch and bound
  - Sankoff-Morel-Cedergren algorithm
  - MALIGN and POY
- Maximum likelihood (ML)
- Bayesian inference (BI)

# 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-invariant sites distributions of rates across sites.

# 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*

# *Phylogenetic Analysis*

- Character state method
  - Maximum parsimony
- Distance method
  - Neighbor-joining and UPGMA method
  - Fitch-Margoliash method
- Maximum likelihood methods
  - determinate evolution model first, then construct system trees

# 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

*Bootstrap*
*Substitution Model*
*Tree Construction*

Evaluate phylogenetic tree

# *General Rule for Method Selection*



(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

# *Make Phylip Packages into Automatic Flow*

# Inside of POWER

# POWER: PhylOgenetic WEb Repeater

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





*Nucl. Acids Res. 2005 33: W553-W556*

# 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.*

# *Service Usage of POWER* *from 2005 July.*



*Near 7,400 Visits*

*More than 136,000 sequences*

# *Automatic On-Line Demonstration*



[http://www.nhri.org.tw/nhri_org/bs/biostat/power.swf](http://www.nhri.org.tw/nhri_org/bs/biostat/power.swf)

# BMC Bioinformatics

Research article

## Linear array of conserved sequence motifs to discriminate protein subfamilies: study on pyridine nucleotide-disulfide reductases

César L Avila[1], Viviana A Rapisarda[1], Ricardo N Farías[1], Javier De Las Rivas[2] and Rosana Chehín*[1]

Address: [1]Departamento Bioquímica de la Nutrición, Instituto Superior de Investigaciones Biológicas (CONICET-UNT) and Instituto de Química Biológica Dr Bernabé Bloj, Chacabuco 461 (4000), San Miguel de Tucumán, Tucumán, Argentina and [2]Instituto de Biología Molecular y Celular ..., Spain

...qf.unt.edu.ar; Ricardo N Farías - rfarias@conicet.gov.ar;

33. Zhang Y, Jock S, Geider K: **Genes of Erwinia amylovora involved in yellow color formation and release of a low-molecular-weight compound during growth in the presence of copper ions.** *Mol Gen Genet* 2000, **264**:233-240.

34. Thompson JD, Gibson TJ, Plewniak F, Jeanmougin F, Higgins DG: **The CLUSTALX windows interface: flexible strategies for multiple sequence alignment aided by quality analysis tools.** *Nucleic Acids Res* 1997, **25**:4876-4882.

35. Eddy SR: **HMMER: Profile hidden Markov models for biological sequence analysis.** 2001 [http://hmmer.wustl.edu/].

36. Page RDM: **TREEVIEW: An application to display phylogenetic trees on personal computers.** *Computer Applications in the Biosciences* 1996, **12**:357-358.

37. **PHYLIP package on POWER** [http://power.nhri.org.tw]

38. Gattiker A, Gasteiger E, Bairoch A: **ScanProsite: a reference implementation of a PROSITE scanning tool.** *Applied Bioinformatics* 2002, **1**:107-108.

39. Crooks GE, Hon G, Chandonia JM, Brenner SE: **WebLogo: A sequence logo generator.** *Genome Research* 2004, **14**:1188-1190.



PHYLIP package on POWER   [http://power.nhri.org.tw]

# *POWER Listed in Bioinfo Portals*

- 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 Tree of life - Collection of WWW pages on phylogeny and biodiversity of organisms
- The PhylOgenetic Web Repeater (POWER) - perform phylogenetic analysis
- NCBI Taxonomy Browser
- NEWT - UniProt Taxonomy Browser

- CluSTr - Automatic classification of UniProtKB proteins into groups of related proteins
- ProtoNet - Classification of the proteins into hierarchical clusters

# 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

*Bootstrap*
*Substitution Model*
*Tree Construction*

Evaluate phylogenetic tree

# *Phylogenetic Analysis*

- Character state method
  - Maximum parsimony
- Distance method
  - Neighbor-joining and UPGMA method
  - Fitch-Margoliash method
- Maximum likelihood methods
  - determinate evolution model first, then construct system trees

# *Flowchart of Analysis*



(Mount, *Bioinformatics*)

# *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 <span style="color:red">maximum likelihood allows additional statistical flexibility</span> 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)

$$L(\tau, ө) = Prob(Data | \tau, ө)$$

Or

Prob(Aligned Sequences | tree, model of evolution)

# *Maximum Likelihood Estimates (MLE)*

- The maximum likelihood estimates (MLE) of $\tau$, $\theta$ are those making the LH function as large as possible

$$\tau, \theta = \max L(\tau, \theta)$$

- Hence, what we usually call the likelihood of the tree is **not the likelihood of the tree**, but **the probability of the data given that the tree is the true tree**.

# *Basic Substitution Model*

- **The models in the GTR family are distinguished by their degree of parameterization**

  **I. Nucleotide frequencies** : $\pi A = \pi C = \pi G = \pi T = 0.25$ ó $\pi A \neq \pi C \neq \pi G \neq \pi T$
  - models assuming = frequencies: JC69; K2P, K3P ...
  - models accomodating ≠ frequencies: F81, HKY85, TrN93, GTR ...

  **II. Substitution rates and types: transitions (ti) and transversions (tv)**



- There are 4 ti and 8 tv substitution types; when **ti/tv ≠ 0.5 there is a substitution rate bias in the data set. Generally ti >> tv.**

- The nucleotide substitution models in the GTR family are also distinguished by the number of rate parameters they use to accomodate the possible substitutions:

| no. rates | model(s) |
|-----------|----------|
| 1 | **JC69** (ti=tv) |
| 2 | **K2P** (ti ≠tv) |
| 3 | **TrN** ó **K3P** (2 ti, 1 tv) |
| 6 | **GTR** (each its own rate) |

# Illustration of DNA substitution Model

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



GTR (for four characters, as is often the case in phylogenetics) requires 6 substitution rate parameters ($x_1 \sim x_6$), as well as 4 equilibrium base frequency parameters.

# *Illustration of Models for DNA*

# *Background*

- Model fitting in phylogenetics has been suggested for many years, yet <span style="color:red">many authors still arbitrarily choose their models</span>, often using the default models implemented in standard computer programs for phylogenetic estimation.

- Here, we want to show the way that a best-fit model can be readily identified. Consequently, given the relevance of models, model fitting should be routine in any phylogenetic analysis that uses models of evolution.

# *Motivation I*

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

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

➢ PALM uses an open-source LAPP (Linux, Apache, PostgreSql, PHP) structure and

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

# *Motivation II*

➢ Model can be selected by following methods including hierarchical likelihood ratio tests (hLRTs), Akaike information criterion (AIC), and Bayesian information criterion (BIC)

➢ PALM can help user to construct the tree with bootstrap based on best substitution model chosen by maximum likelihood.

➢ 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

# *Component Programs of PALM*

➢ PhyML 3.0

➢ ModelTest 3.7

➢ ProtTest 1.4

➢ ClustalW 2.0.3

➢ Seqret (EMBOSS)

# *Models Used in PALM*

- For DNA (24 models)

  - JC69, K80, F81, HKY, TrN, GTR
  - +I, +G

- For Protein (96 models), **Time consuming**

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

# PalmMonitor for Protein Models

96 Models

6 parallel processes

8    8    8    8    8    8

# *Decreasing Time by PALMmonitor*

➢ According the algorithm used in PALM, some models will take a lot of time to calculate the value of maximum likelihood.

- JTT MtREV      3h:00:50
- MtMam MtArt      3h:29:04
- Dayhoff WAG      2h:50:16
- RtREV CpREV      2h:50:19
- Blosum62 VT      2h:49:17
- HIVb HIVw      2h:56:38

➢ All Models      7h:32:10

*Source: 25 sequences with 5000 residues for each*

# *Input and Output of PALM*

- Input format (Protein and DNA)
  - Fasta format
  - Phylip format: Aligned Sequences
  - User tree (if submitted and valid)
- Output
  - Tree topology by php and GD library
  - Tree file in Newick format
  - Aligned Sequence in phylip format
  - Best model selector by PALM

# Flowchart of PALM

# Result of PALM

**PALM**
Phylogenetic reconstruction by Automatic Likelihood Model selector

Home  Demo  Help  Contact

## PALM Result

### Input parameters

| Job ID | 20080525234728289 | Number of Substitution Rate Category | 4 |
|---|---|---|---|
| Sequence Type | DNA | Model Selection Criterion | AICc |
| Number of Bootstrap | 100 | Optimization of Tree Topology | Yes |
| Job Note | test in 100 BS | Optimization of Branch Length | Yes |
| Starting Tree | BIOJ | | |



### Result Information

| Best Model Selected | GTR+G |
|---|---|
| Model Selection Criterion | AICc |
| -lnL | 1903.3464 |
| Number of Estimated Parameters (K) | 9 |

| Model | -lnL | K | AICc | delta | weight | cumWeight |
|---|---|---|---|---|---|---|
| GTR+G | 1903.3464 | 9 | 3825.0603 | 3819.0361 | 0.00e+00 | 1.0000 |
| GTR | 1905.8035 | 8 | 3827.9001 | 3821.8760 | 0.00e+00 | 1.0000 |
| GTR+I+G | 1904.1582 | 10 | 3828.7664 | 3822.7422 | 0.00e+00 | 1.0000 |
| GTR+I | 1905.8112 | 9 | 3829.9897 | 3823.9656 | 0.00e+00 | 1.0000 |
| TrN+G | 1910.6451 | 6 | 3833.4607 | 3827.4365 | 0.00e+00 | 1.0000 |
| HKY | 1912.7156 | 4 | 3833.5120 | 3827.4878 | 0.00e+00 | 1.0000 |
| TrN | 1911.7296 | 5 | 3833.5806 | 3827.5564 | 0.00e+00 | 1.0000 |
| HKY+G | 1911.9691 | 5 | 3834.0596 | 3828.0354 | 0.00e+00 | 1.0000 |
| TrN+I+G | 1910.6479 | 7 | 3835.5234 | 3829.4993 | 0.00e+00 | 1.0000 |
| HKY+I | 1912.7211 | 5 | 3835.5635 | 3829.5393 | 0.00e+00 | 1.0000 |
| TrN+I | 1911.7354 | 6 | 3835.6411 | 3829.6169 | 0.00e+00 | 1.0000 |
| HKY+I+G | 1911.9722 | 6 | 3836.1147 | 3830.0906 | 0.00e+00 | 1.0000 |
| F81+G | 1941.3434 | 4 | 3890.7676 | 3884.7434 | 0.00e+00 | 1.0000 |
| K80 | 1945.1681 | 1 | 3892.3442 | 3886.3201 | 0.00e+00 | 1.0000 |
| F81+I+G | 1941.3442 | 5 | 3892.8098 | 3886.7856 | 0.00e+00 | 1.0000 |
| F81 | 1943.7166 | 3 | 3893.4814 | 3887.4573 | 0.00e+00 | 1.0000 |
| K80+G | 1944.9779 | 2 | 3893.9800 | 3887.9558 | 0.00e+00 | 1.0000 |

### Download Area

| Original File | 20080525234728289 |
|---|---|
| Phylip File | 20080525234728289.phy |
| Phylogenetic Tree (Newick) | tree20080525234728289.txt |
| Statistic data | 20080525234728289_phyml_stat.txt |
| Modelselection Information | Modeltest20080525234728289.out |
| Bootstrap Tree | 20080525234728289_phyml_boot_trees.txt |
| Bootstrap Statistic data | 20080525234728289_phyml_boot_stats.txt |

# *Demonstration of PALM*



Access : http://palm.iis.sinica.edu.tw

# *Bootstrap (BS) Analysis*

- Bootstrap analysis is the most often used 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 due to the original bias connot be corrected by the re-sampling process.

# Input Sequences Make the Tree Different

HIV

# *Future Plans for PALM*

- Integrate more substitution models into PALM
- Improve and optimize the performance of whole pipeline
- MrBayes will be implemented into this system for Bayesian inference.

# Acknowledgement



Daniel, Sheng-Yao Su

Tengi, Huang

Pao-Han Kuo

Chen-Ren Lo

# Protein Network

- ✓ *hp*-DPI
- ✓ *fly*DPI
- ✓ Reconstruction of Human protein network
- ✓ Topological analysis by Hubba

# *Motivations*

- Combine accumulated fragmentary information (experimental interactions) into a systems-level picture (embrace experimental and putative interactions) with spatiotemporal scenarios

- Construct entire network including those interactions can't be done due to experiment limitation (ie. Toxic and membrane proteins can not be tested in Y2H ).

- Understand host and pathogen networks, how they merge during infection

- Provide a multilayered and integrated view to control diseases ranging pathogenic infection to cancer

# *Deciphering Protein into Domains*

- Using the protein-protein interaction (PPI) data set to infer domain-domain interaction (DDI) for specific organism. Then using the predicted DDI set can infer other probable PPI set

- Deciphering the domain interaction will allows us to discover novel interactions between proteins that contain domains with known binding partner.



Protein Interactions



Domain Interactions

# Previous Work I:
## Helicobacter pylori- Database of Protein Interactome

**Letters to Nature**

*Nature* **409**, 211-215 (11 January 2001) | doi: 10.1038/35051615

The protein−protein interaction map of *Helicobacter pylori*

Jean-Christophe Rain[1], Luc Selig[1], Hilde De Reuse[2], Véronique Battaglia[1], Céline Reverdy[1], Stéphane Simon[1], Gerlinde Lenzen[1], Fabien Petel[1], Jérôme Wojcik[1], Vincent Schächter[1], Y. Chemama[1], Agnès Labigne[2] and Pierre Legrain[1]

*Over 1,200 interactions were identified between H. pylori (strain 26695) proteins, connecting 46.6% of the proteome.*

Experimental interactions    Predicted interactions

## The network of whole proteome

© 2004 Division of Biostatistics and Bioinformatics

# *Previous Work, hp-DPI*



**This website can be accessed at** *http://dpi.nhri.org.tw/hp/*

# *Search Result of hp-DPI*

# *Edges Patterns for Interaction*



parent locus: uvrA
        orf: HP0705

child locus:
        orf: HP0677

probability: 1.0
exp. interaction: N

**Solid red Line**

**Inferred data (Prob =1)**

**but not Exp data**

parent locus: copA
        orf: HP1072

child locus: mod
        orf: HP0593

probability: 0.5
real interaction: N

**Dotted Line**

**Inferred data only**

parent locus: rpoB
        orf: HP1198
child locus: trcF
        orf: HP1541
probability: 1.0
exp. interaction: Y

**Solid Line**

**Exp. Data and
inferred  Prob =1**

parent locus: spoT
        orf: HP0775
child locus: spoT
        orf: HP0775
probability: 1.0
exp. interaction: Y

**Solid curve -Self Interaction**

**Exp. Data and inferred Prob =1**

parent locus:
        orf: HP0677
child locus: uvrB
        orf: HP1114
probability: 1.0
exp. interaction: N

**Solid purple Line**

**( related with Level 2 object)**

**Inferred data (Prob =1)**

**but not Exp data**

parent locus: uvrA
        orf: HP0705
child locus:
        orf: HP0452
probability: 0.5
exp. interaction: Y

**Dotted Line**

**Exp data**

**but Inferred data**

**With Prob<1**

# Discover New Research Targets with hp-DPI

Acidic Gastric Juice

pH=2

pH=4    mucus gel layer

pH=7    Hpylori    HCO₃⁻    NH₄⁺

mucus cells

## Urease

$$C=O(NH_2)2 + H^+ + 2H_2O \xrightarrow{urease} HCO_3^- + 2\ NH_4^+$$

*bases*

# *Network of Urease Complex*

# Annotated Protein Function by Interacting Partners



HMG1-box containing protein

DNA binding protein

ATPase component of chromatin remodeling complex

transcription related protein

ATPase component of chromatin remodeling complex

Transcription regulation, chromatin binding

Transcription factor III, Tau

# *The Evolution of the Flagellum*



Structure of a bacterial flagellum. The illustration on the left is a rotationally averaged reconstruction of electron micrographs of purified hook-basal bodies. The names of the various parts are listed in the illustration to the right.

*Reference: Uetz, et al., J. Bacteriol. 2006*

# *Flagellum of H. pylori*



```
orf : HP1157
locus : omp26
description : outer membrane protein
GO :
KEGG : Non-enzymes.
level : 3
```

```
orf : HP0601
locus : flaA
description : flagellin A
GO : structural molecule activity.
      ciliary/flagellar motility.
KEGG : Flagellar assembly.
level : 2
```

HP1542

omp26

HP1154

fliS

flaB

HP1377

cag26   omp9

```
orf : HP0753
locus : fliS
description : flagellar protein
GO : flagella biogenesis.
KEGG : Flagellar assembly.
level : 3
```

```
orf : HP0317
locus : omp9
description : outer membrane protein
GO :
KEGG : Non-enzymes.
level : 3
```

*Reference: Uetz, et al., J. Bacteriol. 2006*

# hp-DPI (http://dpi.nhri.org.tw/hp/)



BIOINFORMATICS

Institution: National Health Research Intitutes  Sign In as Personal Subscriber

Author:  Keyword(s):
Year:    Vol:    Page:

## hp-DPI: *Helicobacter pylori* database of protein interactomes- embracing experimental and inferred interactions

Chung-Yen Lin [1*], Chia-Ling Chen [1], Chi-Shiang Cho [1], Li-Ming Wang [1], Chia-Ming Chang [1], Pao-Yang Chen [1], Chen-Zen Lo [1], and Chao A. Hsiung [1]

[1] Division of Biostatistics and Bioinformatics, National Health Research Institutes. #128, Sec. 2 Yaun-Chio-Yun Rd. Taipei 115, Taiwan

* To whom correspondance should be addressed.
Chung-Yen Lin, E-mail: cylin@nhri.org.tw

# *hp-DPI Selected into*
# *2006 The Molecular Biology Database Collection by NAR*

# Visit Statistics for hp-DPI
## from 2004/11/22 ~ 2008/07/1

# Previous Work II:
# Fly Database of Protein Interactomes



Query in Full Text

## http://flydpi.nhri.og.tw

Search Result &
Statistical Estimation

Network Visualization
& popup annotation

# New Features in FlyDPI



Ping-Pong Search

Full-text Search

Chromosome Location

Gene Categories form GO

Spatiotemporal Scenarios

# Search Results of FlyDPI

## General Search



A snap of the experimental and inferred visualized interaction networks of *D. melanogaster* interactome under specific spatiotemporal scenarios.

## Ping-Pong Search



Map of proteins potentially involved in apoptosis generated by ping-pong search. By the click on the nodes or lines between two query proteins, the advanced option will remove the paths related or confine the paths with the selected nodes or lines

# *Interaction Network Amid Gro And Its Partners*

# FlyDPI --
# (http://flydpi.nhri.org.tw)

## Fly-DPI: database of protein interactomes for *D. melanogaster* in the approach of systems biology

Chung-Yen Lin* [1,2,3] ✉, Shu-Hwa Chen* [4] ✉, Chi-Shiang Cho[1] ✉, Chia-Ling Chen[1] ✉, Fan-Kai Lin[1] ✉, Chieh-Hua Lin[1] ✉, Pao-Yang Chen[1] ✉, Chen-Zen Lo[1] ✉ and Chao A Hsiung[1] ✉

[1]Division of Biostatistics and Bioinformatics, National Health Research Institutes. No. 35 Keyan Rd. Zhunan, Miaoli County 350, Taiwan
[2]Institute of Information Science, Academia Sinica, No. 128 Yan-Chiu-Yuan Rd., Sec. 2, Taipei 115, Taiwan
[3]Institute of Fishery Science, National Taiwan University, No. 1, Sec 4, Roosevelt Road, Taipei, 10617, Taiwan
[4]Stem Cell/Regenerative Medicine Program, Genomics Research Center, Academia Sinica., No. 128 Yan-Chiu-Yuan Rd., Sec. 2, Taipei 115, Taiwan

**BMC Bioinformatics**

# *Visits of FlyDPI* *(Dec 2006- Aug 2008)*

# *Framework for Database of Protein Interactome, DPI*

# Candida albicans Protein Network

- *Candida albicans* is both a <span style="color:red">commensal</span> and <span style="color:red">pathogen</span> of humans that can infect a broad range of body sites.
- Endogenous *C. albicans* infections are established by cells that normally colonise mucosal surfaces or skin as harmless commensals, and that are triggered to cause infection by changes in the host immune system or microflora.



Candida albicans - Corn meal agar, MMBC-UTMB

# C. Albicans with Antifungal Drug Resistance

➢ Question emerging after abusing Antifungal drugs

➢ Identification of novel drug targets by network biology is way to solve the problem.

# *Inferred Protein Interactions by hidden DDIs from Yeast*



Yeast Protein interactions
(Y2H only)

Domain-interaction Network of yeast

*Physiological Scenarios*

Putative *C. albicans* Protein interaction

# CaPTION-
# C. albicans Protein interaction Network

# *Interface of CaPTION (available soon)*

# *Ecoli-DPI*

# Network Biology:
# Hub/ Essential Proteins Identification

# *The Ways to Detect Hubs*

- Degree (Jeong H. *et al*.,2001)
- Bottle Neck (Przulj N. *et al*., 2003)
- Percolation Based **(Vi)** (Chin *et al*., 2003)
- Subgraph centrality **(SC)** (Ernesto, E *et al*., 2005)
- Maximum Connected Component from Neighborhood Induced Subgraph ( **MNCIS** ) (Our team, 2007)
- Maximum Connected Component from Neighborhood Induced Subgraph with Density **(MNCISD)** (Our team, 2007)

a clique with 5 nodes

# *Hub Object Analyzer: Hubba*



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

To appear on NAR 2008 Web issue

# The Relationship of Top 10 in Yeast Complex Network *(PPI from DIP, 2007 Jan)*



Fragile motif in whole network

# Identification Target Proteins and Hubs for Novel Cancer Therapies



Putative Protein Network in Human Cancer

Hub protein can be treated with RNA inference to perturb the network, then stop the progress of tumors

# Inferred Protein Interactions by Conserved and hidden DDIs

Human evolutionary
conserved DDIs

Human experimental
hidden DDIs
(inferred from recent
publications and public DBs )

Physiological
Scenarios

Putative Human
Protein interaction

Pathogen
Network

# *Human Protein Network*

Methodology article

Highly accessed | Open Access

## Reconstruction of human protein interolog network using evolutionary conserved network

**Tao-Wei Huang**[1] ✉, **Chung-Yen Lin**[2,3,4] ✉ and **Cheng-Yan Kao**[1,5] ✉

[1]Department of Computer Science and Information Engineering, National Taiwan University, Taipei 106, Taiwan
[2]Institute of Information Science, Academia Sinica, Taipei 115, Taiwan
[3]Division of Biostatistics and Bioinformatics, National Health Research Institutes, Taipei 115, Taiwan
[4]Institute of Fishery Science, National Taiwan University, Taipei 106, Taiwan
[5]Institute for Information Industry, Taipei 106, Taiwan

✉ author email ✉ corresponding author email

BMC Bioinformatics

$$CS = w_I * \frac{I}{I_R} + w_D * \frac{D}{D_K} + w_T * \frac{T}{T_K} + w_L * \frac{L}{L_K} + w_P * \frac{P}{P_K}$$

PPIs in reference organisms, e.g. mouse, yeast, *etc.*

Interolog using *IP* and *C* scores

PPIs with feature scores ($I, D, T, L, P$) in human

# Confidence score (CS)

- Interolog score (I)   $I_{ij} = w_{ec} * \min(IP_{A_i}, IP_{B_j}) * C_{ab}$

- Domain-domain combination score (D)

$$D = \sum_{j=1}^{2^m-1} \sum_{i=1}^{2^m-1} \frac{N'(pd_i, pd_j)}{N(pd_i, pd_j)} \text{ if } pd_i \in PD_d, pd_j \in PD_d \qquad T = \sum_{i=1}^{79} 1 \text{ if } \log_2 \frac{eA_i}{eA} \geq 1 \text{ and } \log_2 \frac{eB_i}{eB} \geq 1$$

- Tissue specificity score (T)

- Sub-cellular localization score (L)

- Cell-cycle stage score (P)

$$CS = w_I * \frac{I}{I_R} + w_D * \frac{D}{D_K} + w_T * \frac{T}{T_K} + w_L * \frac{L}{L_K} + w_P * \frac{P}{P_K}$$

# *Interactome among Pathogens and Host*

There are 148 nodes and 172 edges in your network. The clustering coeffcient of this network is 0 , and the average path length of this network is 3.24812 .



Source: EBV and Human, Dyer *et al.*, 2008

Source: Our own with Vidal *et al.*, 2007

# *Interactome of Yersinia pestis and Human host*

There are 56 nodes and 49 edges in your network. The clustering coeffcient of this network is 0 , and the average path length of this network is 1.15649 .



http://hub.iis.sinica.edu.tw/Hubba/result.php?ID=upload/2008_04_30_02_36_24

# Protein-protein Interactions For The Ataxia Network

There are 3607 nodes and 6972 edges in your network. The clustering coefficient of this network is 0.0570338 , and the average path length of this network is 4.18696 .



Top 100

Top 10

http://hub.iis.sinica.edu.tw/Hubba/result.php?ID=upload/2008_04_20_00_44_35

# A Protein Interaction Network Links GIT1, an Enhancer of Huntingtin Aggregation, to Huntington's Disease

There are 182 nodes and 592 edges in your network. The clustering coeffcient of this network is 0.23954 , and the average path length of this network is 2.85459 .



Top 90

Top 10

http://hub.iis.sinica.edu.tw/Hubba/result.php?ID=upload/2008_04_16_23_12_26

# Extract Sub-network with Targeted Genes



**Targeted Genes**
ACE
ACE2
AGTR2
AGT
BDKRB2
GHRL
NTS
AGTR1
ARRB2

Relationships among these proteins

Human Interactome, HPRD, 2007/9/01

# *Ongoing Projects*

➢ *Human Stem cell* research and *regenerative medicine* for stemness on expression profile and TF regulatory network (Collaborated with GRC, Academia Sinica)

➢ Protein interactions in the approaches of network analysis and systems biology for human and several model organisms on various spatiotemporal scenarios (granted by NRPGM)

➢ Electronic Lab Notebook (ELN)

# *Electronic Lab Notebook (ELN)*

- Digitalization of Lab notebook from text, gif, raw data, even animations with functions of full text search and security
- Two kinds of version will be provided in the end of this year.
  - For group use: Linux-based version
  - For personal use: USB-ELN, windows/ Mac-based version

# MyBLAST (Customized BLAST Framework)

*http://mybioweb.nhri.org.tw/myblast*

# *MOLAS*

- MicroArray On Line Analysis System (MOLAS): a web-based customizable bioinformatics package designed for manager and analyze massive array data

**New Experiment Design**

**Upload raw data**

*Web interface*     *Exp – Experiment Data*     *Data Normalization*     *Experiment Result and Analysis*

# *Annotated-Feature Extractor from GenBank*



*Coming Soon*

# *Bioinformatics Core for Genomic Medicine and Biotechnology Development*



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

# Selected Publications (2006 - 8)

1) **Lin, C. Y.\*,** Chin, C. H., Wu, H. H., Chen, S. H., Ho, C. W.,\* Ko, M. T.\*, "Hubba: Hub Objects Analyzer : A Framework of Interactome Hubs Identification for Network Biology," Nucleic Acids Res., volume 36, number 2008 Web application Issue, July 2008, *Nucleic Acids Research* Advance Access published online on May 24, 2008 (http://hub.iis.sinica.edu.tw) (SCI/6.945) .

2) Chen, S.H., Lo, C.Z., Tsai, M. C., Hsiung C.A., **Lin, C.Y\*.,** **2008**. "Unique Probe Selector (UPS): A Comprehensive Web Service for Probe Design and Oligo Nucleotide Arrays," To Appear in *BMC Bioinformatics*, (URL: http://array.iis.sinica.edu.tw/ups) (SCI/3.49) .

3) Huang, T. W., **Lin, C. Y\*.,** Kao, C. Y. **2007**. Reconstruction of Human Protein Interolog Network using Evolutionary Conserved Network. *BMC Bioinformatics*. 8:152 (SCI/3.49) .

4) **Lin, C.Y. \*,** Chen S. H., Cho C. S., Chen C. L., Lin F. K., Lin C. H., Chen P. Y., Lo C. Z., and Hsiung C.A., **2006**, "Fly-DPI: Database of Protein Interactomes for *D. melanogaster* in the Approach of Systems Biology.," *BMC Bioinformatics*, 7(5):S18, (SCI/3.49) (URL:http://flydpi.nhri.org.tw)

5) Jiang S. S., Chang I. S., Huang L. W., Chen P. C., Wen C. C., Liu S. C., Chien L. C., **Lin C. Y**., Hsiung C. A., Juang J. L., **2006**"Temporal Transcription Program of Recombinant *Autographa californica* Multiple Nucleopolyhedrosis Virus.," *J. Virol*., 80: 8989-8999. (SCI/ 5.178)

6) Wen, C. C., Wu, Y. J., Huang, Y. H., Chen, W. C., Liu, S. C., Jiang, S. S., Juang, J. L., **Lin, C. Y**., Fang, W. T., Hsiung, C. A., Chang, I. S. **2006**. A Bayes Regression Approach to Array-CGH Data. *Statistical Applications in Genetics and Molecular Biology*. 5(1): art3, (http://www.bepress.com/sagmb/vol5/iss1/art3/), (Medline Index)

7) Chang, C. C., Lin, H. C., Lin, I. P., Chang, T. Y., Chen, H. H., Chen, W. H. Cheng, C. H., **Lin, C. Y**., Liu, S. M. Chang, C. C. Chaw, S. M. **2006**. The Chloroplast Genome of *Phalaenopsis aphrodite* (Orchidaceae): Comparative Analysis of Evolutionary Rate with That of Grasses and Its Phylogenetic Implications. *Mol. Bio. Evol*. 23: 279 - 291 (SCI/ 6.355)

8) Pan W. H., Lynn K. S., Chen C. H., Wu Y. L., Lin C. Y., Chang H.Y. **2006**. Using endophenotypes for pathway clusters to map complex disease genes. *Gen. Epi.* 30(2): 143-154. (SCI/5.42 )

# Acknowledgement

National Health Research Institutes

Chia-Ling Chen
Fan-Kai Lin
Chieh-Hwa Lin
Chia-Ming Chang
Chi-Shiang Cho
Chen-Zen Lo
Yung-Shiang Hwang
Pao-Yang Chen
Ming-Hsin Tasi
Char-Lin Pan
Chao A. Hsiung

Institute of Information Science Academia Sinica

Hsin-Hung Wu
Daniel, Sheng-Yao, Su
Pan-Han, Kuo
Tengi, Huang
Yu-Bin, Wang
Ming-Ta Ko

Department of Computer Science and Information Engineering National Central University

Chia-Hau Kim
Chin-Wen Ho

臺大資訊

Chen-Yen Kao
Tao-Wei Huang

*Thanks for your Attention*