# Protein Families

João C. Setubal

University of São Paulo

Agosto 2012

# Motivation

- *Phytophthora Science* paper [Tyler et al., 2006]

- …Comparison of the [*P. sojae* and *P. ramorum*] genomes reveals a rapid expansion and diversification of <span style="color:red">many protein families</span> associated with plant infection such as <span style="color:red">hydrolases</span>, <span style="color:red">ABC transporters</span>, <span style="color:red">protein toxins</span>, <span style="color:red">proteinase inhibitors</span>, and, in particular, a <span style="color:red">superfamily</span> of 700 proteins with similarity to known oömycete <span style="color:red">avirulence genes</span>.

# The concept of family

- A group where members share one or more characteristics
- In biology: usually descent or function
- Examples

# The family *felidae*

# Oomycota

- Kingdom: Chromalveolata
- Phylum: Heterokontophyta
- Class: **Oomycota**
- Orders (& families)
- Lagenidiales
  - Lagenidiaceae
  - Olpidiosidaceae
  - Sirolpidiaceae
- Leptomitales
  - Leptomitaceae
- Peronosporales
  - Albuginaceae
  - Peronosporaceae
  - Pythiaceae ← Phytophthora
- Rhipidiales
  - Rhipidaceae
- Saprolegniales
  - Ectrogellaceae
  - Haliphthoraceae
  - Leptolegniellaceae
  - Saprolegniaceae
- Thraustochytriales

# Family by function

- The "family" of all effector proteins found in the phytophthora genus

- Members in general do not share an ancestor

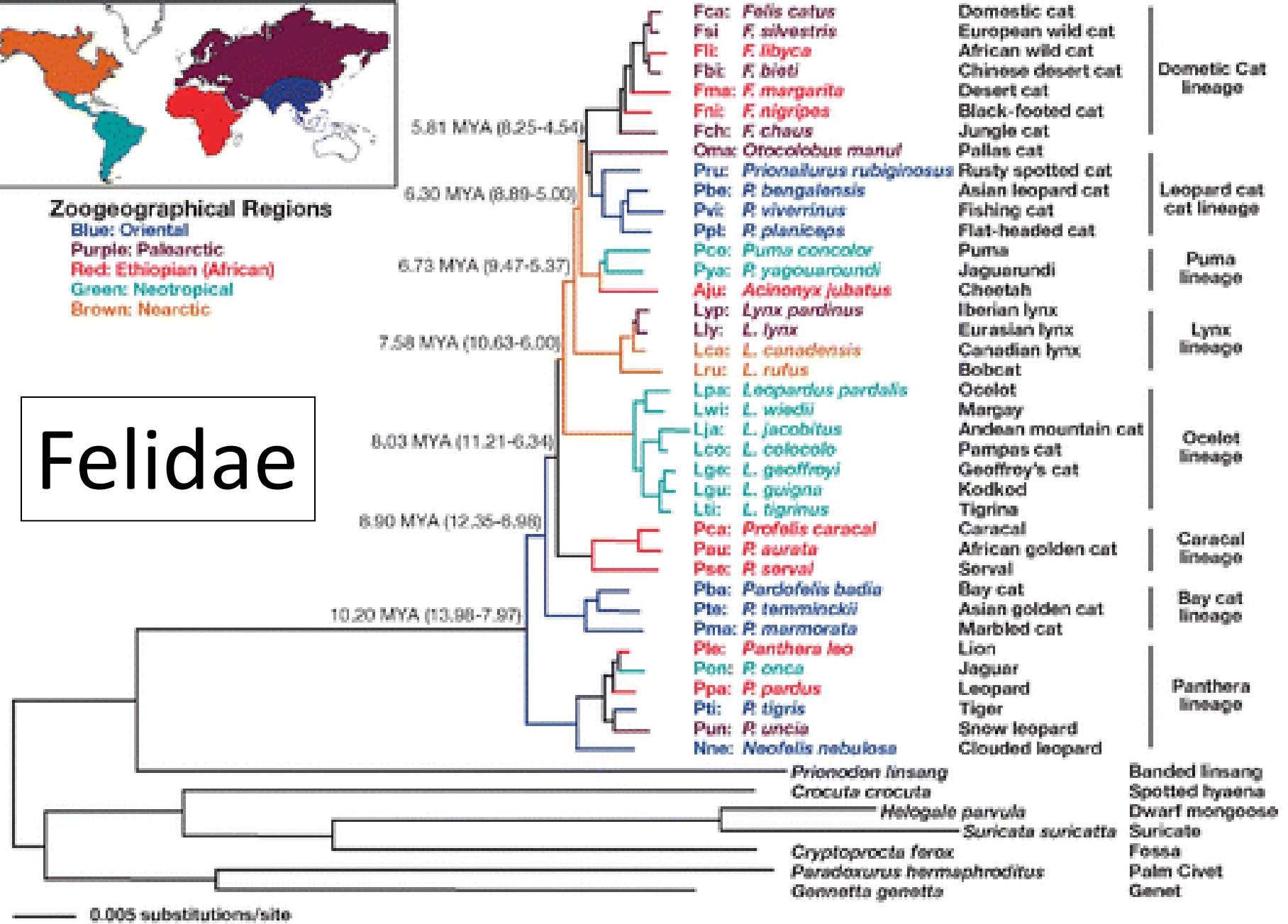- Function is ill-defined in this case

# Protein families

- Shared characteristic: common ancestor
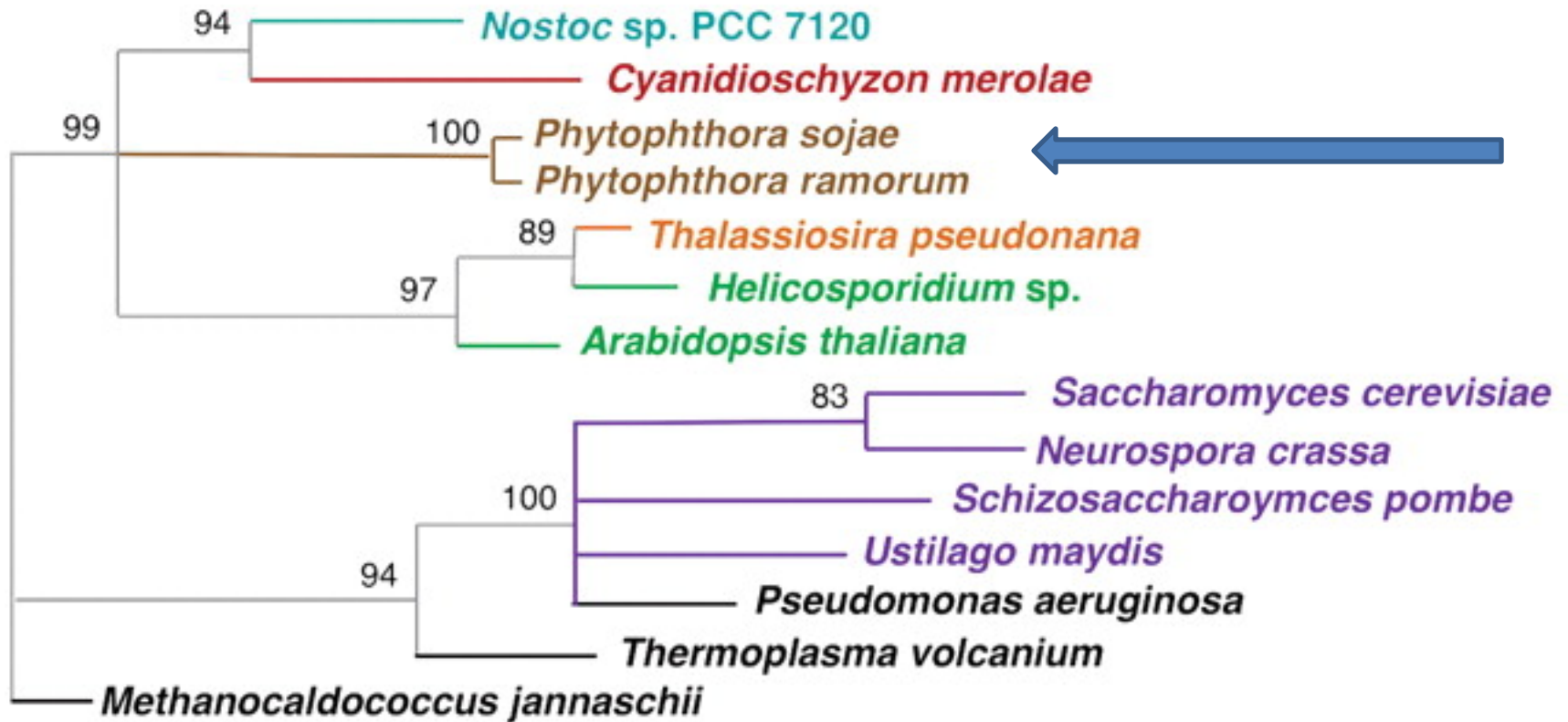- The concept of phylogeny is closely associated with that of family

Felidae

5.81 MYA (8.25-4.54)
6.30 MYA (8.89-5.00)
6.73 MYA (9.47-5.37)
7.58 MYA (10.63-6.00)
8.03 MYA (11.21-6.34)
8.90 MYA (12.35-6.98)
10.20 MYA (13.88-7.97)

| Code | Species | Common name | Lineage |
|---|---|---|---|
| Fca: | Felis catus | Domestic cat | Domestic Cat lineage |
| Fsi: | F. silvestris | European wild cat | |
| Fli: | F. libyca | African wild cat | |
| Fbi: | F. bieti | Chinese desert cat | |
| Fma: | F. margarita | Desert cat | |
| Fni: | F. nigripes | Black-footed cat | |
| Fch: | F. chaus | Jungle cat | |
| Oma: | Otocolobus manul | Pallas cat | |
| Pru: | Prionailurus rubiginosus | Rusty spotted cat | Leopard cat cat lineage |
| Pbe: | P. bengalensis | Asian leopard cat | |
| Pvi: | P. viverrinus | Fishing cat | |
| Ppl: | P. planiceps | Flat-headed cat | |
| Pco: | Puma concolor | Puma | Puma lineage |
| Pya: | P. yagouaroundi | Jaguarundi | |
| Aju: | Acinonyx jubatus | Cheetah | |
| Lyp: | Lynx pardinus | Iberian lynx | Lynx lineage |
| Lly: | L. lynx | Eurasian lynx | |
| Lca: | L. canadensis | Canadian lynx | |
| Lru: | L. rufus | Bobcat | |
| Lpa: | Leopardus pardalis | Ocelot | Ocelot lineage |
| Lwi: | L. wiedii | Margay | |
| Lja: | L. jacobitus | Andean mountain cat | |
| Lco: | L. colocolo | Pampas cat | |
| Lge: | L. geoffroyi | Geoffroy's cat | |
| Lgu: | L. guigna | Kodkod | |
| Lti: | L. tigrinus | Tigrina | |
| Pca: | Profelis caracal | Caracal | Caracal lineage |
| Pau: | P. aurata | African golden cat | |
| Pse: | P. serval | Serval | |
| Pba: | Pardofelis badia | Bay cat | Bay cat lineage |
| Pte: | P. temminckii | Asian golden cat | |
| Pma: | P. marmorata | Marbled cat | |
| Ple: | Panthera leo | Lion | Panthera lineage |
| Pon: | P. onca | Jaguar | |
| Ppa: | P. pardus | Leopard | |
| Pti: | P. tigris | Tiger | |
| Pun: | P. uncia | Snow leopard | |
| Nne: | Neofelis nebulosa | Clouded leopard | |

| | |
|---|---|
| Prionodon linsang | Banded linsang |
| Crocuta crocuta | Spotted hyaena |
| Helogale parvula | Dwarf mongoose |
| Suricata suricatta | Suricate |
| Cryptoprocta ferox | Fossa |
| Paradoxurus hermaphroditus | Palm Civet |
| Genetta genetta | Genet |

0.005 substitutions/site

O'Brien SJ, Johnson WE. 2005.
Annu. Rev. Genomics Hum. Genet. 6:407–29

# 2-isopropylmalate synthase (leucine biosynthesis)



Tyler et al., *Science*, 2006

**A** NPP1 toxin

Bacterial
Fungal
Oomycete { *P. sojae* / *P. ramorum* / Other }

**B**

Percentage of *P. sojae* Avh Proteins vs. Amino acid Identity with Most Similar *P. ramorum* Protein

Avr1b-1

0-30%   30-50%   50-100%

Tyler et al., *Science*, 2006

# Family (or grouping) by function
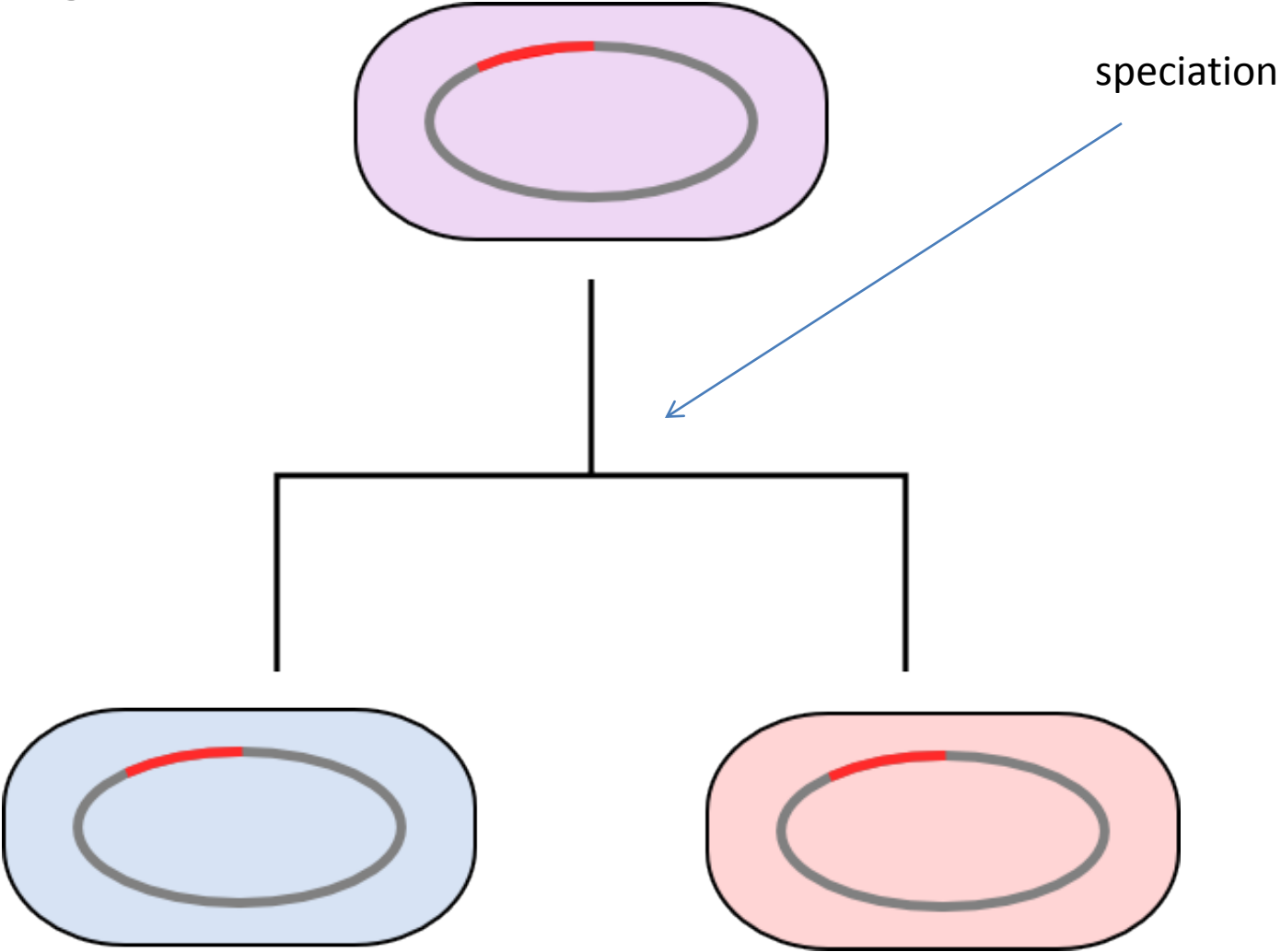
- A controlled vocabulary is necessary
- Gene Ontology

# Protein families by homology

- Aula 1
  - Important concepts
  - How to build a protein family
  - Protein family resources
- Aula 2
  - Phylogeny

# Homology

- Two genes that share an ancestor are said to be homologous

- Often (but wrongly) used with the sense of "similar"

- Similarity does not necessarily mean homology!

- Two kinds of homologous genes
  – Orthologs and paralogs

# Orthologs



speciation

# Paralogs



Figure by C. Lasher

# Protein family

- Operational definitions
  - Two proteins are in the same family if their genes are homologous

- or
  - Two proteins are in the same family if their genes are orthologous

- Will abuse language and mention homologous genes and homologous proteins

# In-paralogs



Figure by C. Lasher

# Homology and function

- We would like orthologous proteins to have the same function

- This is generally but not always the case

- Paralagous genes are more prone to develop new functions
  - Neo-functionalization

- In practice
  - Protein family: homology and shared function
  - Superfamilies and subfamilies

vertebrate class-1

invertebrate

fish class-2B

fish class-2A

vertebrate class-2

0.3

**Phylogenetic tree of the WHAMM proteins**
Kollmar *et al. BMC Research Notes* 2012 **5**:88   doi:10.1186/1756-0500-5-88

# Audiences for this lecture

1.  You have a sequence and you want to build its family

2.  You want to explore and use existing protein family resources for families you are interested in

3.  You have a new genome and you want to place all of its genes into their families

# Audience 1

You have a sequence and you want to build its family

# How to build a family

- Given a protein sequence
  - Determine other members
  - Create multiple alignment
  - Create family signature
  - Create model (Hidden Markov model)

# Pipeline

phylogeny

Input sequence → BLAST → Multiple alignment

NR (NCBI)

curation

HMM

curation

Family signature

Family profile

# Problems & Tools

- BLAST
  - Max e-value and/or minimum identity
  - Minimum coverage
    - 80% query, 80% subject
- PSI-BLAST
- Multiple Alignment
  - Aula 2
- HMM
  - Use Pfam package
- If you don't know what you're doing, don't try this at home! ☺

# Comparação de sequencias

- Similaridade "suficiente" → mesma função
- O que é similaridade?
- O que é "suficiente"?
- Google das sequencias: BLAST
- Basic Local Alignment Search Tool
- Altschul et al., 1990, 1997
- No ano 2000 já tinham mais de 13.000 citações

▸ NCBI/ BLAST/ blastp suite

**Align Sequences Protein BLAST**

| blastn | **blastp** | blastx | tblastn | tblastx |

## Enter Query Sequence

BLASTP programs search protein subjects using a protein query. more..

**Enter accession number(s), gi(s), or FASTA sequence(s)** ⊙          Clear          **Query subrange** ⊙

```
>s
MQLNLAMGAVADGDRAPKACDAACSEAAGDKSAMMHDALFERFSARLKAQVGPEVYASWFA
RLKLHTVSKSVVRFTVPTTFLKSWINNRYMDLITSLVQSEDPDVLKVEILVRSASRPVRPA
QTEERAQPVQEVGAAPRNKSFIPSQSATAPAAQPMAAQATLRQGGSGPLFGSPLDTRFTFD
TFVEGSSNRVALAAAKTIAEAGAGAVRFNPLFIHAGVGLGKTHLLQAIANAAIDSPRNPRV
VYLTAEYFMWRFATAIRDNDALTLKDTLRNIDLLVIDDMQFLQGKMIQHEFCHLLNMLLDS
```

From [          ]

To [          ]

**Or, upload file**      [          ] Browse.. ⊙

**Job Title**      [ s                                                    ]

Enter a descriptive title for your BLAST search ⊙

☑ **Align two or more sequences** ⊙

## Enter Subject Sequence

**Enter accession number, gi, or FASTA sequence** ⊙          Clear          **Subject subrange** ⊙

```
>t
MRSRGISACIQENNYETPETNADARCLETTCEELFKNVSSKLEDQVGSDVYASWFQRLKFR
SVSHNIVYLSVPTNFLKAWIKNRYIDTITKLFQESISSIQGVEIIVRSAAIMPSETSSSSA
IAHTTAKPPIINTGKISTIQGKQSINPVFGSPLDSKFVFSNFIEGPSNRVALAAAHTIAEE
NSSSCTVRFNPLFIHASVGLGKTHLLQAIANAAIKKQNNLRVVYLTAEYFMWRFATAIRDN
YALNFKDCLRNIDLLLIDDMQFLQGKLIQHEFCHLLNSLLDSAKQIVAAADRPPSELESLD
```

From [          ]

To [          ]

**Or, upload file**      [          ] Browse.. ⊙

## Program Selection

**Algorithm**      ⦿ blastp (protein-protein BLAST)

Choose a BLAST algorithm ⊙

```
>lcl|35099 t
Length=499

 Score =  604 bits (1558),  Expect = 0.0, Method: Compositional matrix adjust.
 Identities = 301/499 (60%), Positives = 365/499 (73%), Gaps = 25/499 (5%)

Query  21   DAACSEAAGDKSAMMHDALFERFSARLKAQVGPEVYASWFARLKLHTVSKSVVRFTVPTT  80
            DA C E   ++         LF+  S++L+ QVG +VYASWF RLK  +VS ++V  +VPT
Sbjct  23   DARCLETTCEE-------LFKNVSSKLEDQVGSDVYASWFQRLKFRSVSHNIVYLSVPTN  75

Query  81   FLKSWINNRYMDLITSLVQSEDPDVLKVEILVRSASRPVRPAQTEERAQPVQEVGAAPRN  140
            FLK+WI NRY+D IT L Q     +  VEI+VRSA+  + P++T                +
Sbjct  76   FLKAWIKNRYIDTITKLFQESISSIQGVEIIVRSAA--LMPSETS--------------S  119

Query  141  KSFIPSQSATAPAAQPMAAQATLRQGGSGPLFGSPLDTRFTFDTFVEGSSNRVALAAAKT  200
             S I   +A P           +    P+FGSPLD++F F  F+EG SNRVALAAA T
Sbjct  120  SSAIAHTTAKPPIINTGKISTIQGKQSINPVFGSPLDSKFVFSNFIEGPSNRVALAAAHT  179

Query  201  IAEAGAGA--VRFNPLFIHAGVGLGKTHLLQAIANAAIDSPRNPRVVYLTAEYFMWRFAT  258
            IAE  + +  VRFNPLFIHA VGLGKTHLLQAIANAAI    N RVVYLTAEYFMWRFAT
Sbjct  180  IAEENSSSCTVRFNPLFIHASVGLGKTHLLQAIANAAIKKQNNLRVVYLTAEYFMWRFAT  239

Query  259  AIRDNDALTLKDTLRNIDLLVIDDMQFLQGKMIQHEFCHLLNMLLDSAKQVVVAADRAPW  318
            AIRDN AL  KD LRNIDLL+IDDMQFLQGK+IQHEFCHLLN LLDSAKQ+V AADR P
Sbjct  240  AIRDNYALNFKDCLRNIDLLLIDDMQFLQGKLIQHEFCHLLNSLLDSAKQIVAAADRPPS  299

Query  319  ELESLDPRVRSRLQGGMAIEIEGPDYDMRYEMLNRRMGSARQDDPSFEISDEILTHVAKS  378
            ELESLD R+RSRLQGG+A+ +    D +MR  +L  R+  A++D+P   IS+EIL  VA++
Sbjct  300  ELESLDSRIRSRLQGGVAVPLGAHDIEMRLTILKNRLKMAKKDNPKLYISEEILQRVAQT  359

Query  379  VTASGRELEGAFNQLMFRRSFEPNLSVDRVDELLSHLVGSGEAKRVRIEDIQRIVARHYN  438
            VT SGREL+GAFNQL+FR SFEP L++  VDELLSHLV +GE K++RIEDIQR+V++HYN
Sbjct  360  VTTSGRELDGAFNQLVFRNSFEPVLTIKMVDELLSHLVSAGETKKIRIEDIQRMVSKHYN  419

Query  439  VSRQELVSNRRTRVIVKPRQIAMYLAKMLTPRSFPEIGRRFGGRDHTTVLHAVRKIEDLI  498
            +SR +L+SNRR R IV+PRQIAMYL+K++TPRSFPEIGRRFG RDHTTVLHAVRKIE  +
Sbjct  420  ISRTDLLSNRRVRTIVRPRQIAMYLSKIMTPRSFPEIGRRFGDRDHTTVLHAVRKIEKSM  479

Query  499  SGDTKLGHEVELLKRLINE    517
              DT +  EVELLKRLI+E
Sbjct  480  EKDTVIKKEVELLKRLISE    498
```

Plot of lcl|35097 vs 35099

# Buscando no GenBank

# Lista de hits



**Descriptions**

Legend for links to other resources: **U** UniGene **E** GEO **G** Gene **S** Structure **M** Map Viewer PubChem BioAssay

Sequences producing significant alignments:

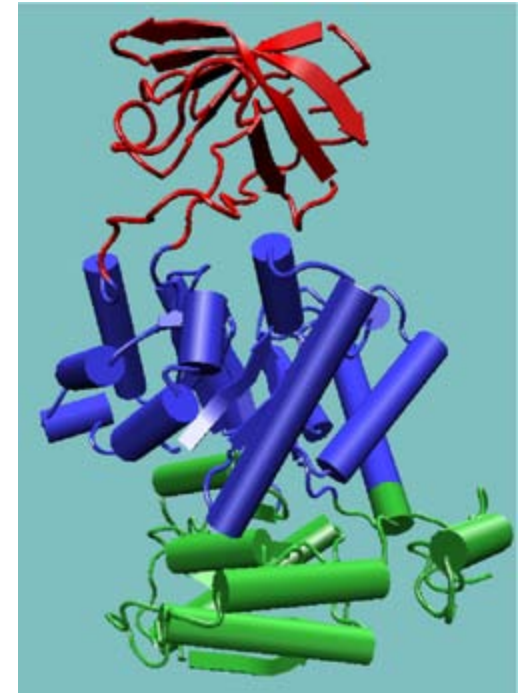| Accession | Description | Max score | Total score | Query coverage | E value | Max ident | Links |
|---|---|---|---|---|---|---|---|
| YP_004062317.1 | chromosome replication initiator DnaA [Candidatus Liberibacter solanacearum | 785 | 785 | 99% | 0.0 | 76% | G |
| YP_003065040.1 | dnaA gene product [Candidatus Liberibacter asiaticus str. psy62] >gb|ACT57 | 755 | 755 | 99% | 0.0 | 76% | G |
| YP_002543179.1 | chromosomal replication initiation protein [Agrobacterium radiobacter K84] >g | 615 | 615 | 97% | 0.0 | 61% | G |
| YP_765982.1 | dnaA gene product [Rhizobium leguminosarum bv. viciae 3841] >sp|Q1MMD6. | 613 | 613 | 93% | 0.0 | 63% | G |
| YP_001976569.1 | chromosomal replication initiation protein [Rhizobium etli CIAT 652] >gb|ACE8 | 612 | 612 | 93% | 0.0 | 63% | G |
| YP_002973852.1 | dnaA gene product [Rhizobium leguminosarum bv. trifolii WSM1325] >gb|ACS. | 612 | 612 | 93% | 0.0 | 63% | G |
| YP_467907.1 | chromosomal replication initiation protein [Rhizobium etli CFN 42] >gb|ABC89| | 611 | 611 | 93% | 0.0 | 63% | G |
| YP_002279530.1 | dnaA gene product [Rhizobium leguminosarum bv. trifolii WSM2304] >gb|ACI5 | 608 | 608 | 93% | 0.0 | 64% | G |
| EGP58677.1 | chromosomal replication initiation protein [Agrobacterium tumefaciens F2] | 607 | 607 | 95% | 0.0 | 61% | |
| EHS51424.1 | Chromosomal replication initiator protein dnaA [Rhizobium sp. PDO1-076] | 605 | 605 | 94% | 0.0 | 62% | |
| EHH08270.1 | chromosomal replication initiation protein [Agrobacterium tumefaciens CCNW( | 600 | 600 | 93% | 0.0 | 62% | |
| YP_002548273.1 | chromosomal replication initiation protein [Agrobacterium vitis S4] >gb|ACM3 | 601 | 601 | 94% | 0.0 | 61% | G |
| NP_353356.2 | chromosomal replication initiation protein [Agrobacterium tumefaciens str. C5 | 598 | 598 | 95% | 0.0 | 60% | G |
| ZP_08526429.1 | chromosomal replication initiation protein [Agrobacterium sp. ATCC 31749] >s | 595 | 595 | 93% | 0.0 | 61% | |
| YP_004277622.1 | chromosome replication initiator DnaA [Agrobacterium sp. H13-3] >gb|ADY63 | 593 | 593 | 93% | 0.0 | 61% | G |
| YP_001325697.1 | dnaA gene product [Sinorhizobium medicae WSM419] >gb|ABR58862.1| chror | 590 | 590 | 93% | 0.0 | 63% | G |
| ZP_02164856.1 | chromosomal replication initiation protein [Hoeflea phototrophica DFL-43] >gb | 578 | 578 | 93% | 0.0 | 61% | |
| ZP_05929413.1 | chromosomal replication initiator protein dnaA [Brucella abortus bv. 3 str. Tul | 578 | 578 | 93% | 0.0 | 60% | |
| P35890.3 | RecName: Full=Chromosomal replication initiator protein DnaA | 573 | 573 | 93% | 0.0 | 62% | |
| NP_384474.1 | chromosomal replication initiation protein [Sinorhizobium meliloti 1021] >ref|Y | 574 | 574 | 93% | 0.0 | 62% | G |
| AAA26258.1 | dnaA [Sinorhizobium meliloti] >gb|AAA91097.1| dnaA [Sinorhizobium meliloti] | 574 | 574 | 93% | 0.0 | 62% | |
| YP_001608612.1 | dnaA gene product [Bartonella tribocorum CIP 105476] >emb|CAK00617.1| c | 573 | 573 | 92% | 0.0 | 59% | G |
| AFL48605.1 | chromosomal replication initiator protein DnaA [Sinorhizobium fredii USDA 257 | 571 | 571 | 93% | 0.0 | 62% | |
| YP_004547390.1 | unnamed protein product [Sinorhizobium meliloti AK83] >gb|AEG51776.1| Chr | 573 | 573 | 93% | 0.0 | 62% | G |
| YP_002824558.1 | chromosomal replication initiation protein [Sinorhizobium fredii NGR234] >gb|A | 571 | 571 | 93% | 0.0 | 62% | G |
| CBI78638.1 | chromosomal replication initiator protein DnaA [Bartonella sp. AR 15-3] | 572 | 572 | 92% | 0.0 | 60% | |
| YP_002971177.1 | chromosomal replication initiator protein DnaA [Bartonella grahamii as4aup] > | 571 | 571 | 92% | 0.0 | 60% | G |
| ZP_10237186.1 | chromosomal replication initiation protein, partial [Nitratireductor aquibiodomi | 569 | 569 | 97% | 0.0 | 56% | |

# Outros tipos de análise

- Alinhamento múltiplo



- Construção de árvore filogenética
- Construção de uma assinatura de uma família
  - Modelo oculto de Markov (HMM)
- Predição de estrutura

# Issues

- Two proteins may share a domain and be unrelated
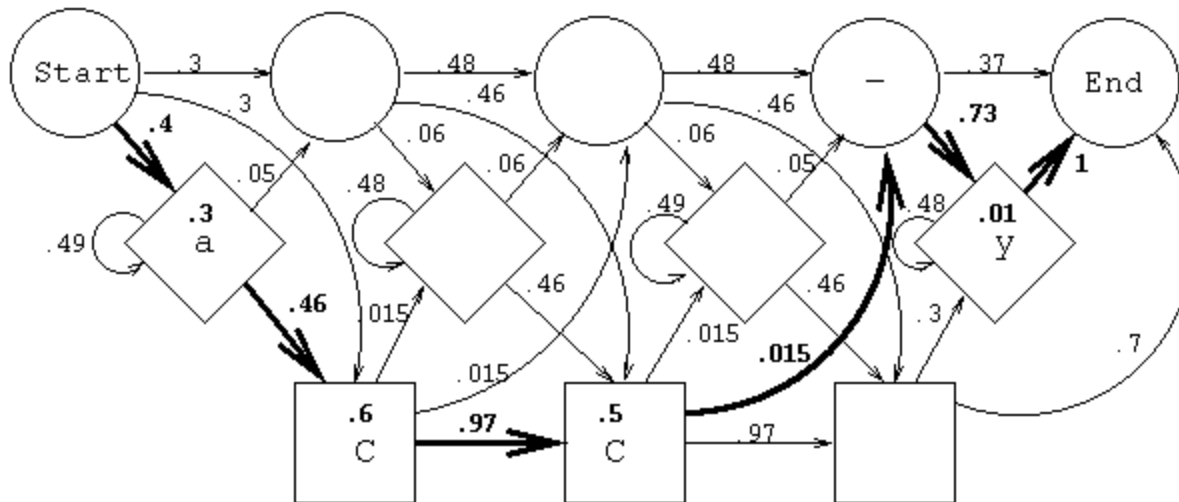  - BLAST false positives
- Multiple domain architectures



Pyruvate kinase

http://en.wikipedia.org/wiki/File:1pkn.png

# HMM

- They capture what is 'essential' to define the family

# Audience 2

You want to explore and use existing protein family resources for families you are interested in

# Protein family resources



Clusters of orthologous groups (COG, KOG, eggNOG)

KEGG orthologs

## Pfam 26.0 (November 2011, 13672 families)

The Pfam database is a large collection of protein families, each represented by **multiple sequence alignments** and **hidden Markov models (HMMs)**. **More...**

keyword search Go

| QUICK LINKS | YOU CAN FIND DATA IN PFAM IN VARIOUS WAYS... |
|---|---|
| **SEQUENCE SEARCH** | Analyze your protein sequence for Pfam matches |
| **VIEW A PFAM FAMILY** | View Pfam family annotation and alignments |
| **VIEW A CLAN** | See groups of related families |
| **VIEW A SEQUENCE** | Look at the domain organisation of a protein sequence |
| **VIEW A STRUCTURE** | Find the domains on a PDB structure |
| **KEYWORD SEARCH** | Query Pfam by keywords |
| **JUMP TO** | enter any accession or ID   Go   Example |

Enter any type of accession or ID to jump to the page for a Pfam family or clan, UniProt sequence, PDB structure, etc.

Or view the help pages for more information

## Recent Pfam blog posts

⊠Hide th

### Does my family of interest have a determined 3D protein structure? (posted 9 May 2012)

Two related questions that we are often asked via the Pfam helpdesk is 'Which families have a known three-dimensional structure?' and 'Why is a particular a PDB structure not found in Pfam'. You may thir that there are obvious answers to these questions – but as with many things in life the answer is not [...]

### TreeFam is back with a new release ! (posted 27 March 2012)

# Query by accession

# Query by sequence search



Sequence search results

Show the detailed description of this results page.

We found **2** Pfam-A matches to your search sequence (**1** significant and **1** insignificant). You did not choose to search for Pfam-B matches.

Show the search options and sequence that you submitted.

Return to the search form to look for Pfam domains on a new sequence.

**Significant Pfam-A Matches**

Show or hide all alignments.

| Family | Description | Entry type | Clan | Envelope | | Alignment | | HMM | | Bit score | E-value | Predicted active sites | Show/hide alignment |
|--------|-------------|------------|------|----------|-----|-----------|-----|------|-----|-----------|---------|----------------------|---------------------|
| | | | | Start | End | Start | End | From | To | | | | |
| Ala_racemase_N | Alanine racemase, N-terminal domain | Domain | CL0036 | 5 | 229 | 6 | 228 | **2** | **216** | 154.3 | 3e-45 | 34 | Show |

**Insignificant Pfam-A Matches**

Show or hide all alignments.

| Family | Description | Entry type | Clan | Envelope | | Alignment | | HMM | | Bit score | E-value | Predicted active sites | Show/hide alignment |
|--------|-------------|------------|------|----------|-----|-----------|-----|------|-----|-----------|---------|----------------------|---------------------|
| | | | | Start | End | Start | End | From | To | | | | |
| ATP-cone | ATP cone domain | Domain | n/a | 96 | 172 | 97 | 160 | **10** | **54** | 10.5 | 0.55 | n/a | Show |

Comments or questions on the site? Send a mail to **pfam-help@sanger.ac.uk**. Our **cookie policy**.
**The Wellcome Trust**

## PANTHER
Classification System

**Quick links**

Whole genome function views

Gene expression tools

cSNP tools

Upload multiple gene IDs

Community Curation

My Workspace

HMM scoring

Downloads

Genome statistics

Site map

**Newsletter subscription**

Enter your Email:

[                    ]

Subscribe

### KEYWORD SEARCH

[ All ▼ ]   [                    ]   [ Go ]

### SEQUENCE SEARCH

Enter a protein sequence: (?)

[                                    ]

Sequence query limits: Protein - 50kb

[ Submit ]

The PANTHER (**P**rotein **AN**alysis **TH**rough **E**volutionary **R**elationships) Classification System is a unique resource that **classifies genes by their functions**, using published scientific experimental evidence and evolutionary relationships to predict function even in the absence of direct experimental evidence. Proteins are **classified by expert biologists** according to:

- Gene families and subfamilies, including annotated phylogenetic trees
- Gene Ontology classes: molecular function, biological process, cellular component
- PANTHER Protein Classes
- Pathways, including diagrams

PANTHER is part of the Gene Ontology Reference Genome Project.

PANTHER is supported by a research grant from the National Institute of General Medical Sciences [grant GM081084] and maintained by the Thomas lab at the University of Southern California.

What can I do on the PANTHER site?
Guide to getting started

**News**

(March 16, 2012)

PANTHER 7.2 is released.

Click for additional info.

**Publications**

How to cite PANTHER

"PANTHER version 7: improved phylogenetic trees, orthologs and collaboration with the Gene Ontology Consortium." Mi, et al.

"Applications for protein sequence-function evolution data: mRNA/protein expression analysis and coding SNP scoring tools." Thomas, et al.

"PANTHER: a library of protein families and subfamilies indexed by function." Thomas, et al.

# PANTHER HMM SEQUENCE SCORING RESULTS ⓘ

The top scoring HMM is reported, along with the E-value (the number of expected false-positive hits expected). If the E-value is less than 1e-3, no hits are reported.

PANTHER Hit: PROLINE SYNTHETASE CO-TRANSCRIBED BACTERIAL HOMOLOG PROTEIN (PTHR10146)
HMM E-value score: 2.2e-113 ●●● ⓘ

```
Sequence  Domain   seq-f seq-t    hmm-f hmm-t       score  E-value
--------  -------  ----- -----    ----- -----       -----  -------
sequence   1/1       1   234 []    11   253 ..     387.5 2.2e-113
```

Alignments of top-scoring domains:
sequence: domain 1 of 1, from 1 to 234: score 387.5, E = 2.2e-113

```
                    *->lgvaanLakVlerikaaaakagRdppavrLvAVSKTkPaelileayd
                       ++va+nL++V++++++a+ak+ R    + +LvAVSKTkP+e+++eay+
      sequence     1   MAVAKNLLAVRAKVAEAVAKSARQ-QQCTLVAVSKTKPVEDLQEAYE 46

                       aGqRhFGENYvQElleKaplLpdlcpdikWHFIGhLQsNKvkkll.gvpn
                       a qRhFGENY+QEl++KaplLp    d+kWH+IGh+QsNK+k+l+++vpn
      sequence    47   ADQRHFGENYIQELVQKAPLLPK---DVKWHYIGHVQSNKAKPLVrDVPN 93

                       ldmvhsvDslklAdklnkaaaklkglgkplkvlvQVNtsGEesKsGvppe
                       l++v++vDs+k+A++lnka  ++  ++++l+v+vQVNts Ee+KsG++ +
      sequence    94   LFVVETVDSIKIANALNKASGEF--RSEKLNVMVQVNTSEEEQKSGIDAD 141

                       ElpellkhvlkkcpnLellGLMTIGpfdgdlekgpnpdFalLaklrkevc
                       +el++h+++ c++L+l GLMTIG++++  ++     F +L+++rk+v+
      sequence   142   GSVELAQHIVSSCEHLNLTGLMTIGRYGDTTSE----CFDRLVACRKRVA 187

                       kklglnpkllELSMGMSgDfelAIeaGsTlVRvGsaIFGeRdypkkp<-*
                       +++g  +  l LSMGMSgDfelAI  GsT+VRvGs+IFG+R+y +k+
      sequence   188   EAIGKAETDLDLSMGMSGDFELAISCGSTHVRVGSTIFGARNYANKE    234
```

# PANTHER Classification System

Home | Browse | Genes and orthologs | **Trees and HMMs** | Pathways | Ontologies | Tools | Workspace

Search Families | Download HMMs |

## Search

PANTHER families ▾

[        ] Go

## Quick links

Whole genome function views

Gene expression tools

cSNP tools

Upload multiple gene IDs

Community Curation

My Workspace

HMM scoring

Downloads

Genome statistics

Site map

## Newsletter subscription

Enter your Email:

[        ]

Subscribe

## PANTHER FAMILY INFORMATION ⓘ

| | |
|---|---|
| Family: | PROLINE SYNTHETASE CO-TRANSCRIBED BACTERIAL HOMOLOG PROTEIN (PTHR10146) |
| Subfamilies: | 1 |
| PANTHER Links: | Tree  MSA |
| GO Molecular Function: | catalytic activity |
| GO Biological Process: | metabolic process ↳ primary metabolic process ↳ cellular amino acid and derivative metabolic process ↳ cellular amino acid metabolic process |
| GO Cellular Component: | |
| PANTHER protein class: | |
| Pathway Categories: | No pathway information available |
| Genes: | 45 |
| HMM Length | 273 |
| Downloads: | HMM (HMMER format) |

## GENES ASSIGNED TO THIS FAMILY

| Species | Count |
|---|---|
| Anopheles gambiae | 1 |
| Aquifex aeolicus vf5 | 1 |
| Arabidopsis thaliana | 2 |

# PANTHER TREE VIEWER ⑦ • close window

Family Name : PROLINE SYNTHETASE CO-TRANSCRIBED BACTERIAL HOMOLOG PROTEIN (PTHR10146)

Tree   MSA

## Tree

| | Grid | MSA |

```
                                                                        125|                    150|
        SF0-PANTR|ENSEMBL=ENSPTRG00000020160|ENSEMBL=ENSPTRP0000003451   DLPAIqpRLVAVSKTKPADMVI.EAYG...YGQRTFGENYVQ..
        SF0-HUMAN|ENSEMBL=ENSG00000147471|UniProtKB=O94903               DLPAIqpRLVAVSKTKPADMVI.EAYG...HGQRTFGENYVQ..
        SF0-MACMU|ENSEMBL=ENSMMUG00000000945|ENSEMBL=ENSMMUP000000001    DLPAIqpRLVAVSKTKPADMVI.EAYG...HGQRTFGENYVQ..
        SF0-RAT|RGD=1308962|NCBI=XP_224947                               LGLPAIqpRLVAVSKTKPTEMVI.EAYG...HGQRTFGENYVQ..
        SF0-MOUSE|MGI=MGI-1891207|UniProtKB=Q9Z2Y8                       DLPAIqpRLVAVSKTKPADMVI.EAYG...HGQRTFGENYVQ..
        SF0-CANFA|ENSEMBL=ENSCAFG00000                                   DLPAIqpRLVAVSKTKPVEMVI.EAYC...HGQRTFGENY---..
        SF0-BOVIN|ENSEMBL=ENSBTAG00000011075|UniProtKB=Q3T0G5            DLPAIqpRLVAVSKTKPADMVI.EAYS...HGQRTFGENYVQ..
        SF0-MONDO|ENSEMBL=ENSMODG00000010863|ENSEMBL=ENSMODP000000136    DLPAIqpRLVAVSKTKPADMVI.EAYA...HGQRSFGENYVQ..
        SF0-ORNAN|ENSEMBL=ENSOANG00000014073|ENSEMBL=ENSOANP0000002219   DLPAVqpRLVAVSKTKPADMVI.EAYI...HGQRSFGENYVQ..
        SF0-CHICK|ENTREZ=426770|NCBI=XP_424381                           IGLPDMqpRLVAVSKTKPAEMVL.DAYS...HGQRSFGENYVQ..
        SF0-DANRE|ENSEMBL=ENSDARG00000060288|ENSEMBL=ENSDARP00000079227  TLPCIppRLVAVSKTKPPEMVV.EAYK...HGQRNFGENYVN..
        SF0-FUGRU|ENSEMBL=ENSTRUG00000003881|ENSEMBL=ENSTRUP00000009144  ALPAV1pRLVAVSKTKPPDLVV.EAYR...QGQRNFGENYVN..
        SF0-CIOIN|ENSEMBL=ENSCING00000016316|ENSEMBL=ENSCINP00000028340  TVPTVqpILVAVSKTKPLSLIK.QAYD...AGQRHFGENYLK..
        SF0-CIOIN|ENSEMBL=ENSCING00000002094|ENSEMBL=ENSCINP00000004266  KLPTVqpILVAVSKTKPLSLIK.QAYD...AGQRHFGENYLK..
        SF0-ANOGA|ENSEMBL=AGAP001621|ENSEMBL=AGAP001621-PA               LSNAPkpLLIAVSKTKPVDLIL.NAYS...VGQRDFGENYVQ..
        SF0-DROME|FB=FBgn0039751|UniProtKB=Q9VA97                        KEVQAArpLLVAVSKTKPAEAVI.EAYE...GGQRDFGENYVQ..
        SF0-CAEBR|ENTREZ=5630972|NCBI=XP_001679808                       ATKRC..RLVAVSKTKSAEMIE.SCFS...QGQRHFGENYVQ..
        SF0-CAEEL|WB=WBGene00017286|UniProtKB=P52057                     ATKRC..RLVAVSKTKSADLIE.ACYS...QNQRHFGENYVQ..
        SF0-NEUCR|ENTREZ=3872168|UniProtKB=Q873K9                        -GRPV..RLVAVSKLKPANDIL.ALHQapqVQHAHFGENYAQ..
        SF0-EMENI|ENTREZ=2870088|UniProtKB=Q5AXG3                        -PKEP..RLVAVSKLKPASDIL.ALHNpp.TAHSHFGENYLQ..
        SF0-ASHGO|ENTREZ=4619930|UniProtKB=Q75B73                        RRSEV..LLLAVSKLKPASDVA.ILYEe..MGLRHFGENYVQ..
        SF0-YEAST|SGD=S000000132|UniProtKB=P38197                        NASKI..LLLVVSKLKPASDIQ.ILYD...HGVREFGENYVQ..
        SF0-SCHPO|GeneDB_Spombe=SPAC644.09|UniProtKB=Q9P6Q1              -GRNV..LLVAVSKFHPVETLM.EAYN...AGQRHFGENYMQ..
        SF0-DICDI|dictyBase=DDB_G0278713|UniProtKB=Q1ZXI6                --HNV..KLVAVSKTKPTEMIR.ILYD...KGHRHFGENYIQ..
        SF0-ENTHI|ENTREZ=3409454|NCBI=XP_655138                          REKPV..CLIAVSKTKPKEAIQ.HLYNv..YNHRVFGENYIQ..
        SF0-ARATH|TAIR=locus=2008910|NCBI=NP_563897                      GSDQI..RVVAVSKTKPVSLIR.QVYD...AGQRSFGENYVQ..
        SF0-ARATH|TAIR=locus=2116387|NCBI=NP_567760                      DAERV..RVLPVSKTKPVSLIR.QIYD...AGHRCFGENYVQ..
        SF0-ORYSJ|ENTREZ=4337823|UniProtKB=Q0DKP7                        APESV..RVVAVSKTKPVGVIR.GVYD...AGHRCFGENYVQ..
        SF0-CHLRE|ENTREZ=5715016|UniProtKB=A8HP79                        -THPV..RLVAVSKTKPAEALQ.EAYD...AGQRVFGENYVQar
        SF0-LEIMA|ENTREZ=5652157|UniProtKB=Q4QAZ0                        -NRRV..TLIAVSKTKSPACLL.NLYN...LGQRVFGENYVQ..
        SF0-TETTHI|ENTREZ=4512506|NCBI=XP_001021721                      -TSDC..TIWCASKTKDLELLQ.QAYD...AGLRHFGENYVD..
```

Berkeley
Phylogenomics
Group

# PhyloFacts 3.0.2

PhyloFacts release PF3.0.2 contains 7,337,238 protein sequences from 99,254 unique taxa (including strains) across 92,800 families (25,446 grouped by PFAM domain and 67,354 grouped by multi-domain architecture agreement). More ...

| | |
|---|---|
| **SEQUENCE ACCESSION SEARCH** | Query PhyloFacts by UniProt accession or identifier |
| **ORTHOLOG IDENTIFICATION** | PhyloFacts Orthology Group: phylogenetic orthologs |
| **JUMP TO PHYLOFACTS FAMILY** | View PhyloFacts family alignments, trees, and annotations |
| **PHYLOFACTS-PFAM SEARCH** | Query PhyloFacts by Pfam accession (PhyloFacts-Pfam Project) |
| **GENOME COVERAGE** | View coverage of key species in PhyloFacts |
| **STATISTICS** | View PhyloFacts coverage statistics |
| **DOWNLOADS** | Download PhyloFacts data |
| **CITING PHYLOFACTS** | How to cite PhyloFacts |

**PhyloFacts statistics**
**7.3M unique proteins across the Tree of Life**

- Bacteria (5.4M, 73.7%)
- Eukaryotes (1.6M, 22.0%)
- Archaea (157k, 2.1%)
- Viruses (153k, 2.1%)
- Unclassified (4.3k, 0.05%)

# Phylofacts
# query by sequence search

**PHOG0274269_00186 – Proline synthetase co-transcribed bacterial protein**

| | |
|---|---|
| PHOG tree: | View tree |
| Pfam domains: | Ala_racemase_N |
| Taxonomic distribution: | stramenopiles |
| PhyloFacts family: | bpg0243724 |
| Alignment: | Global |
| Number of sequences: | 4 |
| Alignment length: | 296 |

## PhyloFacts Orthology Group Members

| Gene ID | Species | Description | Swiss | GO | EC | KEGG | Lit. |
|---------|---------|-------------|-------|----|----|------|------|
| D8LNZ4 | Ectocarpus siliculosus (Brown alga) | Putative uncharacterized protein | | | | | |
| B7GC89 | Phaeodactylum tricornutum (strain CCAP 1055/1) | Predicted protein | | | | | |
| D0MS28 | Phytophthora infestans T30-4 | Proline synthetase co-transcribed bacterial protein | | | | | 🔲 |
| B8BUT1 | Thalassiosira pseudonana (Marine diatom) | Predicted protein | | | | | 🔲 |

Download As CSV

# Resource federation: InterPro

# Not as easy as it may sound…

- Specific protein families may not be consistent across resources
- Most families (MSAs, trees, HMMs) in these resources are not manually curated
  - Domains in Pfam-A are curated
  - TIGRfams are curated
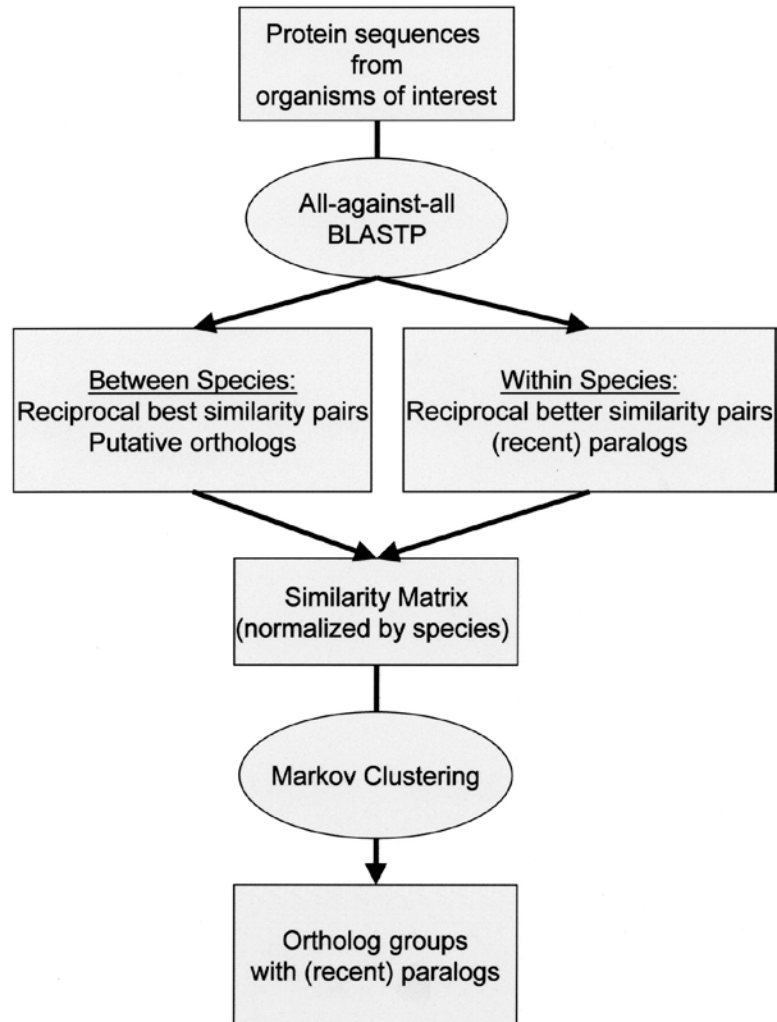  - HAMAP families are curated

# Audience 3

You have a new genome and you want to place all of its genes into their respective families

# Solutions

- Build one at a time (impractical)
- orthoMCL
- multiParanoid

# orthoMCL pipeline



**Li Li et al. Genome Res. 2003; 13: 2178-2189**

# OrthoMCL DB
## Ortholog Groups of Protein Sequences

- Home
- About OrthoMCL ▾
- Data ▾
- Search ▾
- Tools ▾

# Welcome to OrthoMCL DB

Ortholog Groups of Protein Sequences from Multiple Genomes!

Current Release:
Version: 5
Number of Genomes: 150
Number of Protein Sequences: 1398546
Number of Ortholog Groups: 124740

## Search for Groups

- by IDs, Keyword, or PFam domain
- by Phyletic Pattern
- by Phyletic Pattern - Advanced
- by Group Properties
- Query History - Groups

## Search for Sequences

- by IDs, Keyword, Taxonomy or PFam domain
- by BLAST Search
- Query History - Sequences

## Tools

- Assign your proteins to OrthoMCL groups

## OrthoMCL Software

## Database Download

# Final remarks

- The need for experimental results

- Conserved hypotheticals
  - The domino effect