

# Probability Mapping and Bipartition Analysis to Study Genome Histories

J. Peter Gogarten

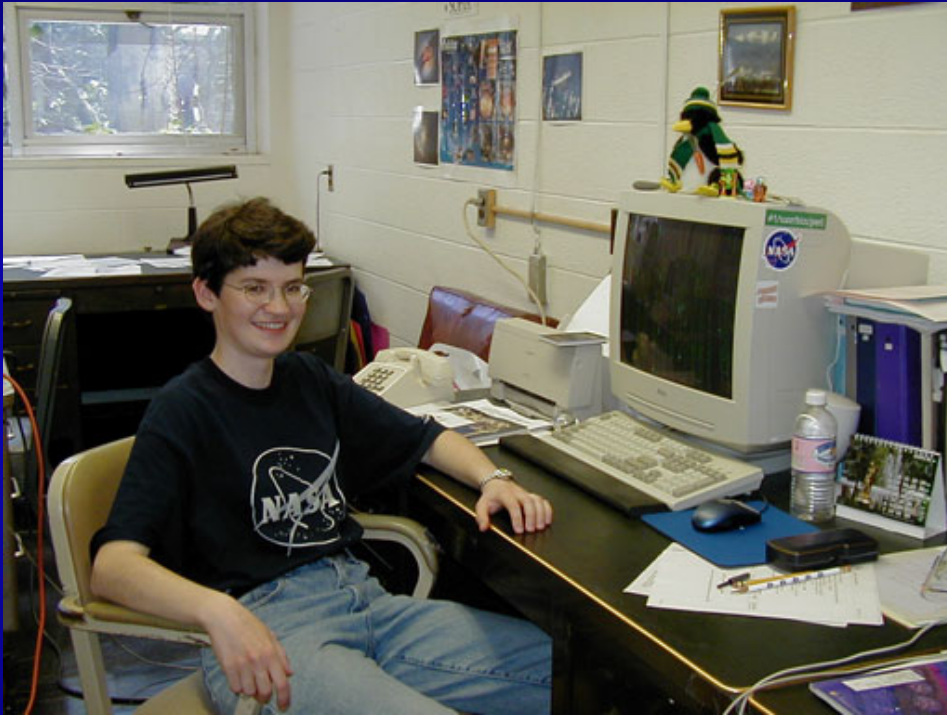
and

Olga Zhaxybayeva

Dept. of Molecular and Cell Biology, Univ. of Connecticut

# Acknowledgements

Olga Zhaxybayeva:



**HGT:**

Lutz Hamel (URI)  
Paul Lewis (UConn)  
Robert Blankenship (ASU)  
Jason Raymond (ASU)  
Ford Doolittle (Dalhousie)  
Jeffery Lawrence (Pittsburgh)  
Gary Olsen (Urbana)

**Coalescence:**

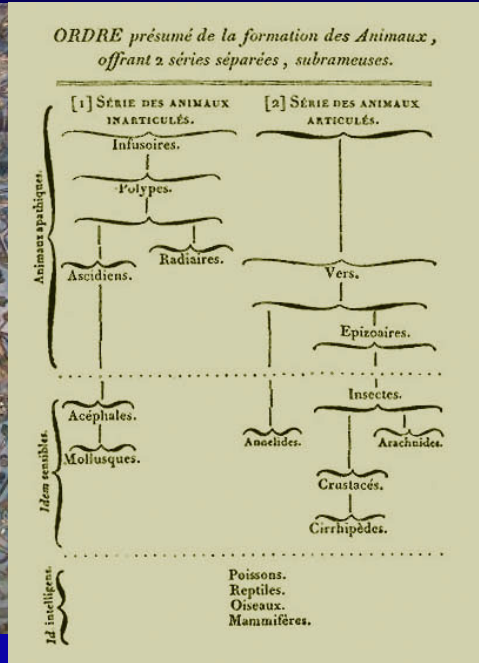
Andrew Martin (U of C Boulder)  
Joe Felsenstein (U of Wash)  
Hyman Hartman (MIT)  
Yuri Wolf (NCBI)

NASA Exobiology Program  
NSF Microbial Genetics

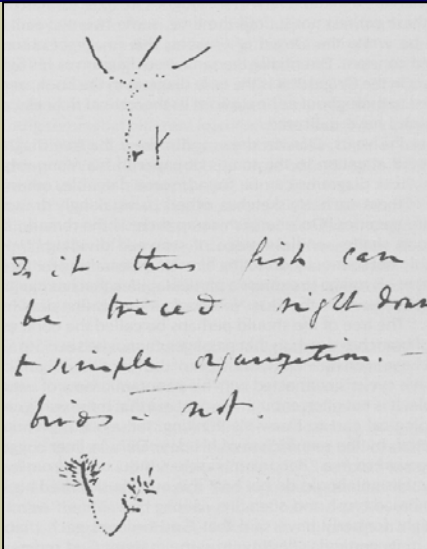
# Trees as a Visualization of Evolution



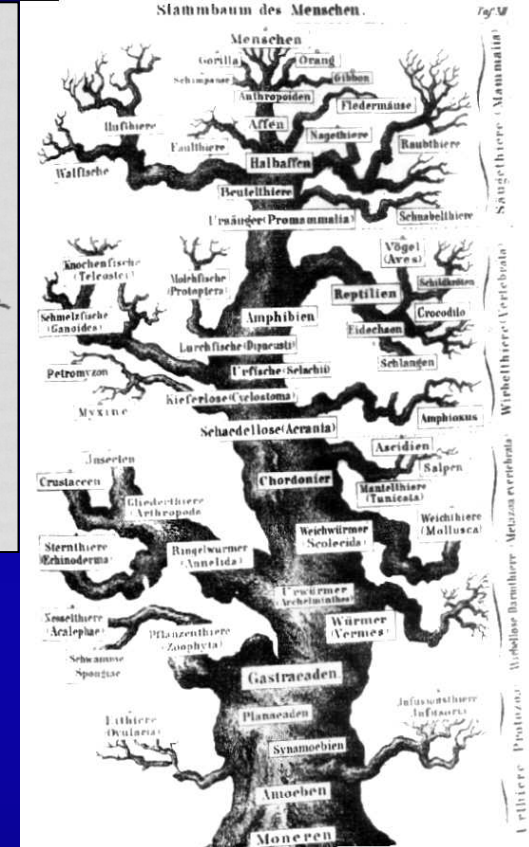
Genealogy  
(Church Ceiling,  
Santo Domingo,  
Oaxaca)



Lamarck's Tree of Life  
(1815)



Page B26 from  
Charles Darwin's  
(1809-1882)  
notebook (1837):



Lebensbaum  
(German for  
"Tree of Life")  
from  
Ernst Haeckel, 1874

# SSU-rRNA Tree of Life

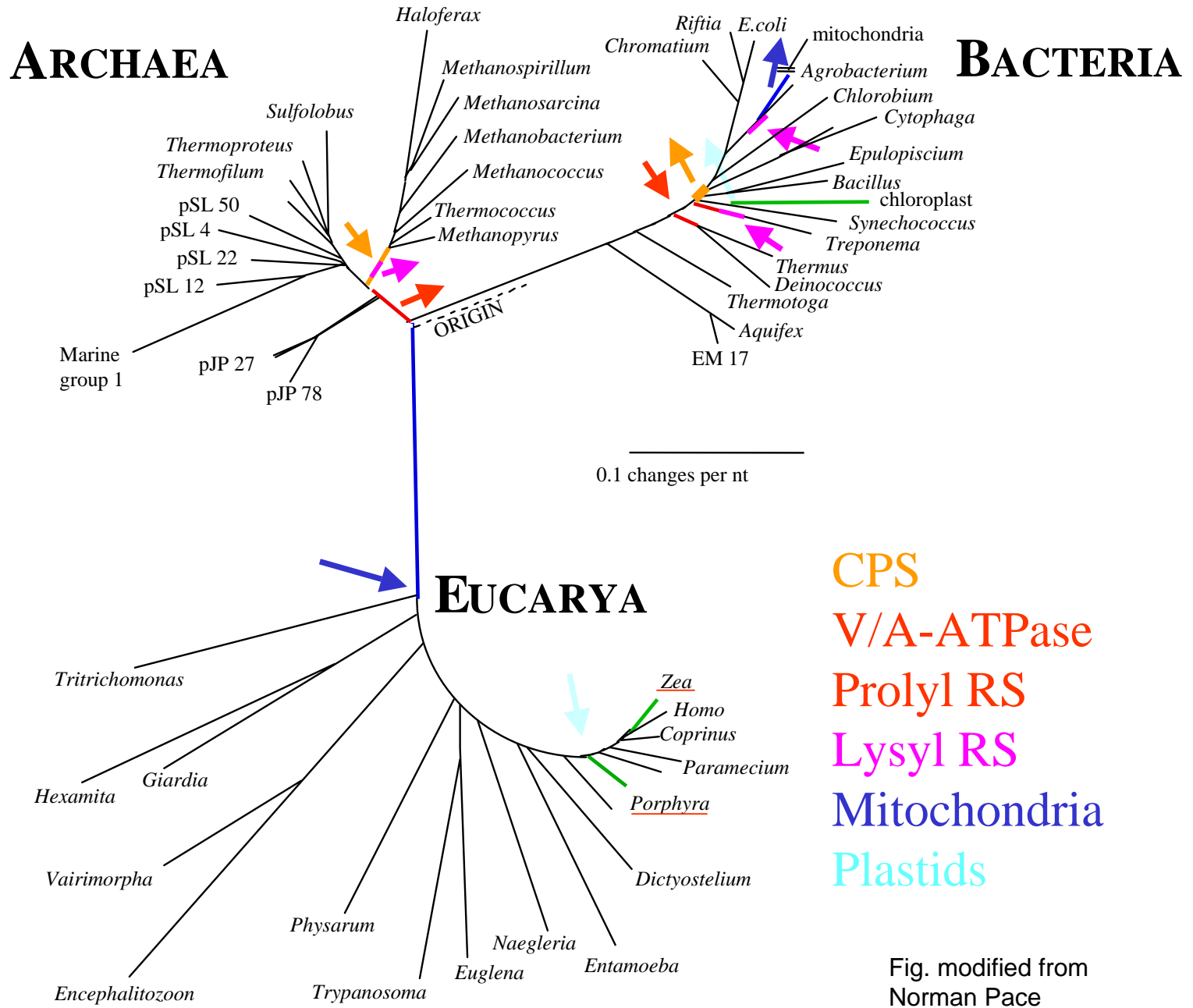


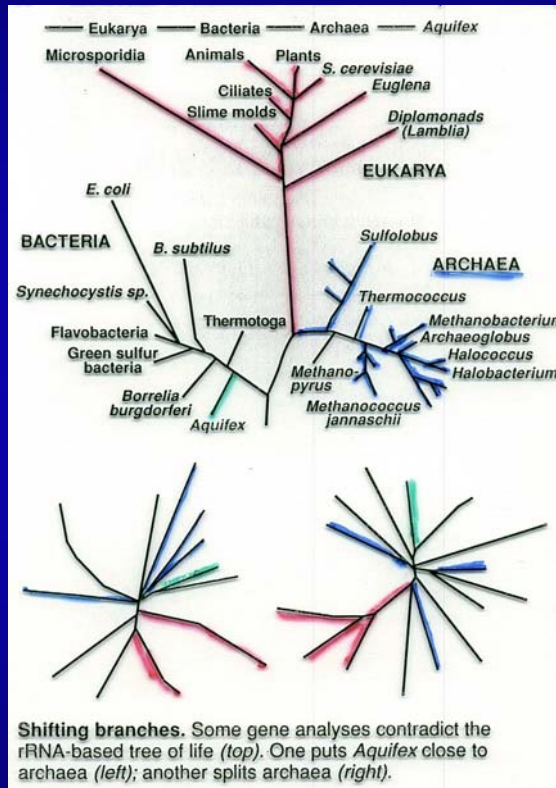
Fig. modified from Norman Pace



RESEARCH NEWS

# Genome Data Shake Tree of Life

New genome sequences are mystifying evolutionary biologists by revealing unexpected connections between microbes thought to have diverged hundreds of millions of years ago



Horizontal Gene Transfer leads to Mosaic Genomes, where different parts of the genome have different histories.

Publicly Available Prokaryotic Genomes:

181 - completed

236 - in progress

Science, **280** p.672ff (1998)

(as of September 8, 2004)

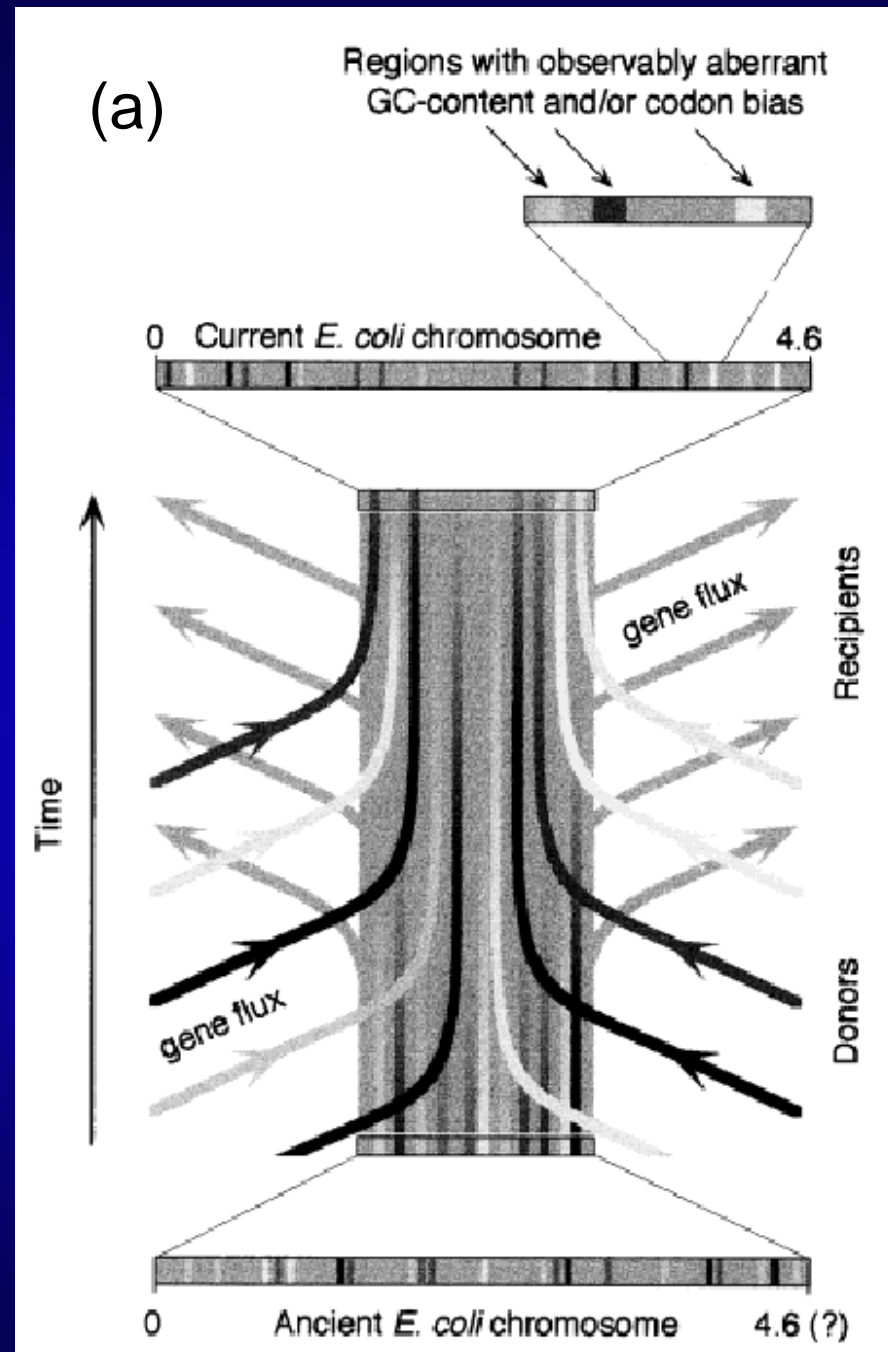
## Transferred genes can be detected using:

(a) unusual composition,

(b) the comparison between closely related species, or

(c) conflicting molecular phylogenies.

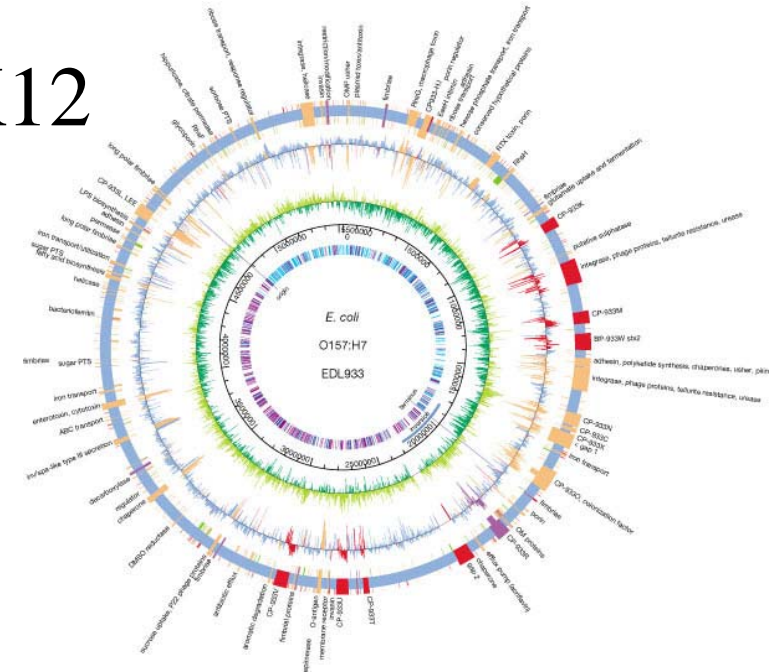
**From Bill Martin  
BioEssays 21 (2), 99-104.**



# *E. coli* O157:H7 versus *E. coli* K12

- divergence about 4.5 million years ago

"We find that lateral gene transfer is far more extensive than previously anticipated. In fact, **1,387 new genes encoded in strain-specific clusters** of diverse sizes were found in O157:H7."



Common: **4,100,000** bp; **3,574** protein-coding genes  
(about 95% identical each on the nucleotide level)

Only in O157:H7: **1,340,000** bp; **1,387** protein-coding genes

Only in K12: **530,000** bp, **528** protein-coding genes

From: **Perna *et al.*** (2001) Nature 409: 529-33

see also **Hayashi *et al.*** (2001) DNA Res. 8:11-22

Welch RA, *et al.*

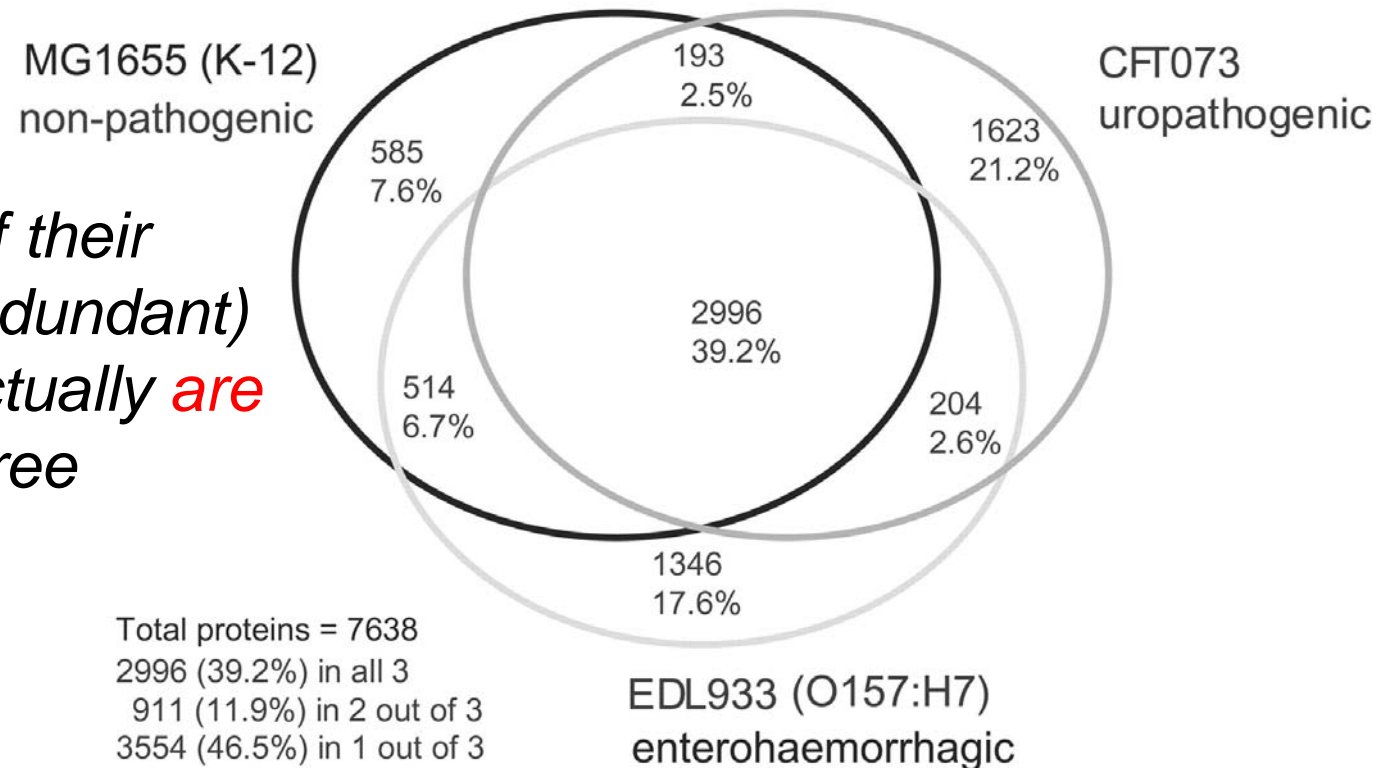
Proc Natl Acad Sci U S A. 2002; 99:17020-4

*Escherichia coli*, strain CFT073, uropathogenic

*Escherichia coli*, strain EDL933, enterohemorrhagic

*Escherichia coli* K12, strain MG1655, laboratory strain,

“... only **39.2%** of their combined (nonredundant) set of proteins actually **are common** to all three strains.”





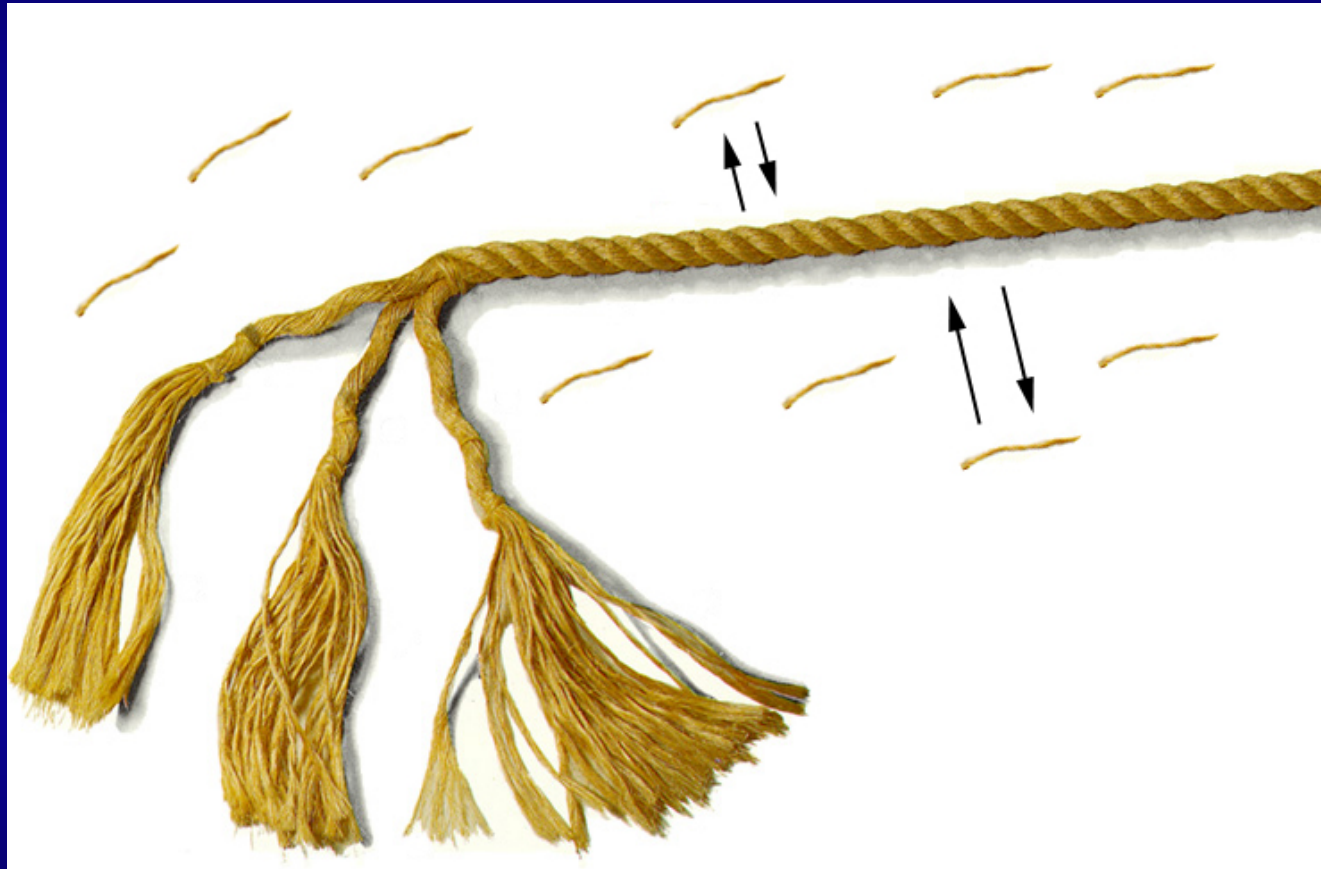
# What is an “organismal lineage” in light of horizontal gene transfer?

Over very **short** time intervals an organismal lineage can be defined as the majority consensus of genes.

This definition only “fails”, if two organisms make co-equal contributions (e.g. endosymbiosis).

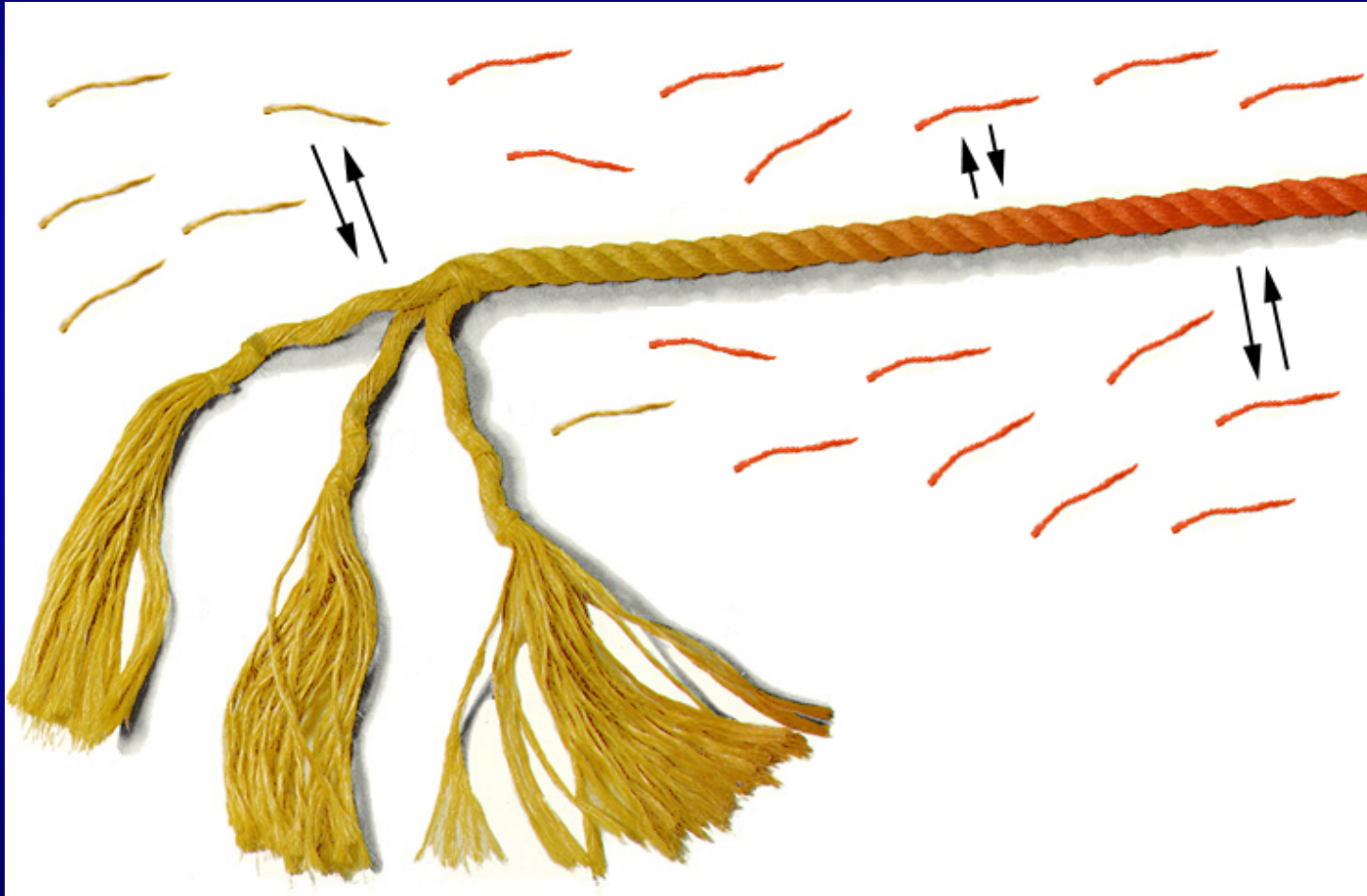
# Rope as a metaphor to describe an organismal lineage (Gary Olsen)

Individual fibers = genes that travel for some time in a lineage.

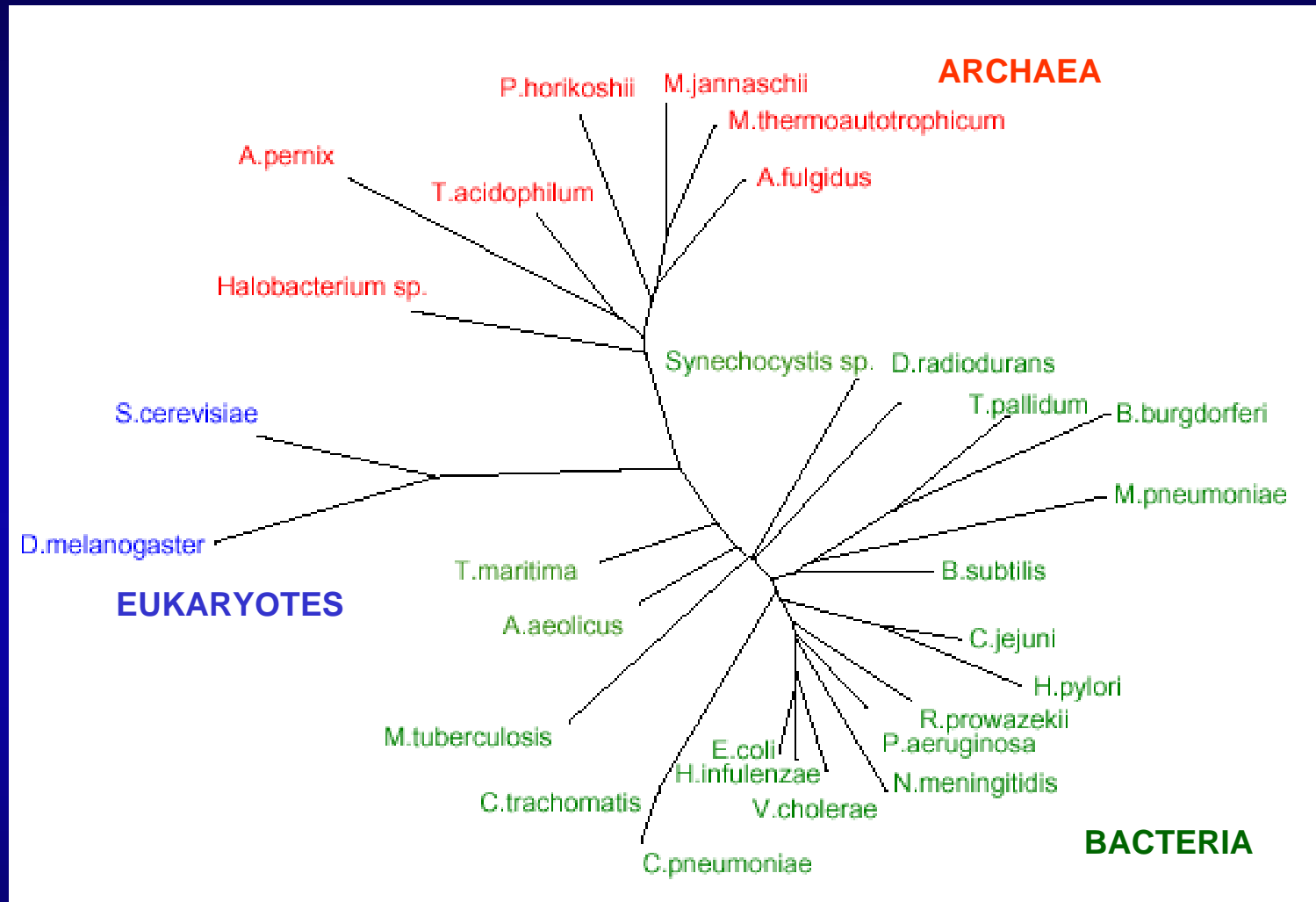


While no individual fiber present at the beginning might be present at the end, the rope (or the organismal lineage) nevertheless has continuity.

However, the genome as a whole will acquire the character of the incoming genes (the rope turns solidly red over time).

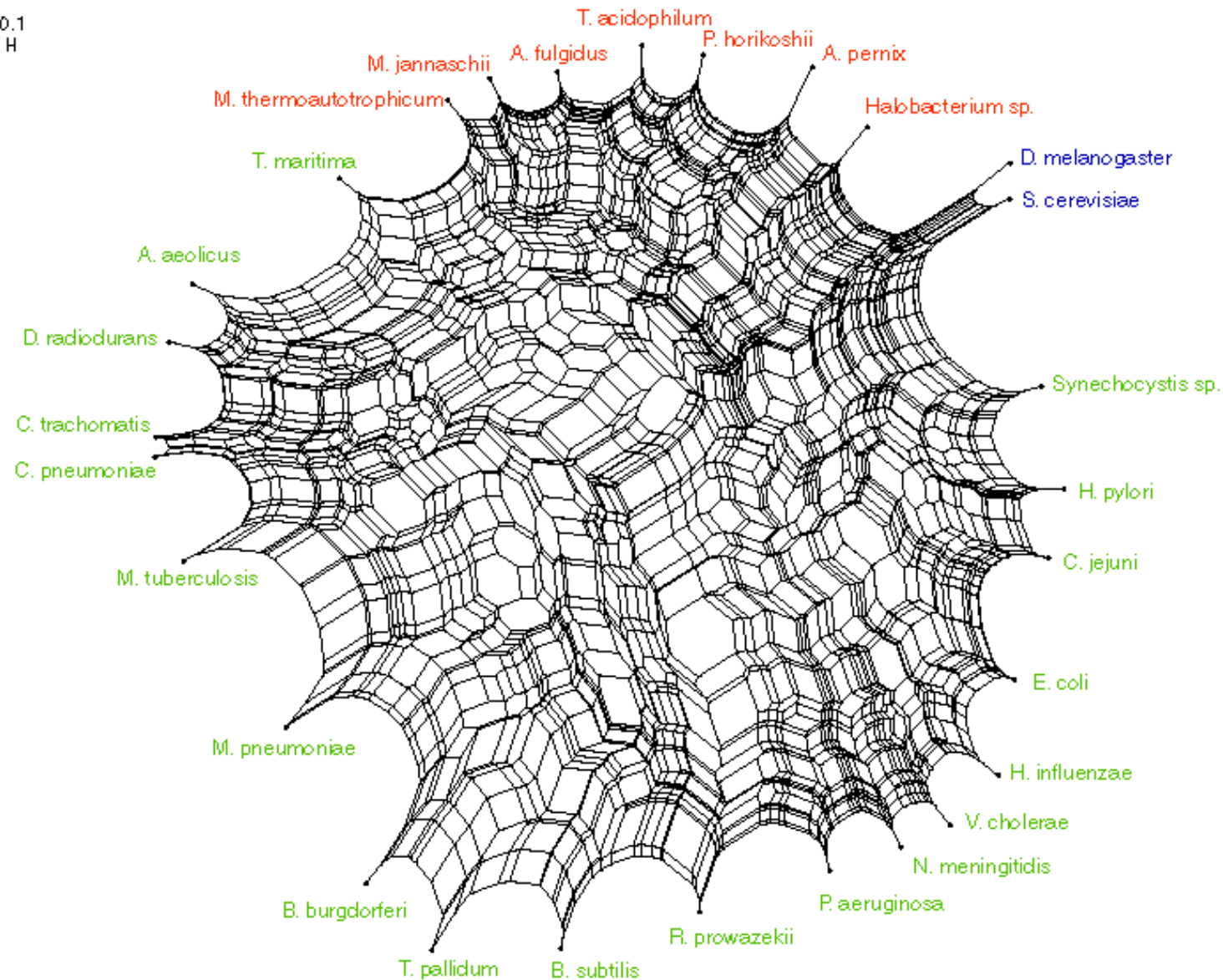


# Genome Content Tree



**Other genome content trees:** Tekaia et al. (1999) *Genome Res* **9**:550-557; Snel et al. (1999) *Nat Genet* **21**:108-110; Lin & Gerstein (2000) *Genome Res* **10**:808-818; Fitz-Gibbon & House (1999) *Nucleic Acids Res* **27**:4218-4222 and (2002) *J Mol Evol* **54**:539-47; Charlebois et al. (2003) *Nature* **421**:217; Wolf et al. (2001), *BMC Evol. Biol* **1**:8

0.1  
H



Same data as before, but network calculated using NeighborNet (David Bryant 2002, <http://www.mcb.mcgill.ca/~bryant/NeighborNet/>)



# Visualization of Mosaic Genome Content

# Bayes' Theorem



Reverend Thomas Bayes  
(1702-1761)

Likelihood

describes how well the model predicts the data

$$P(\text{model}|\text{data}, I) = P(\text{model}, I) \frac{P(\text{data}|\text{model}, I)}{P(\text{data}, I)}$$

Posterior  
Probability

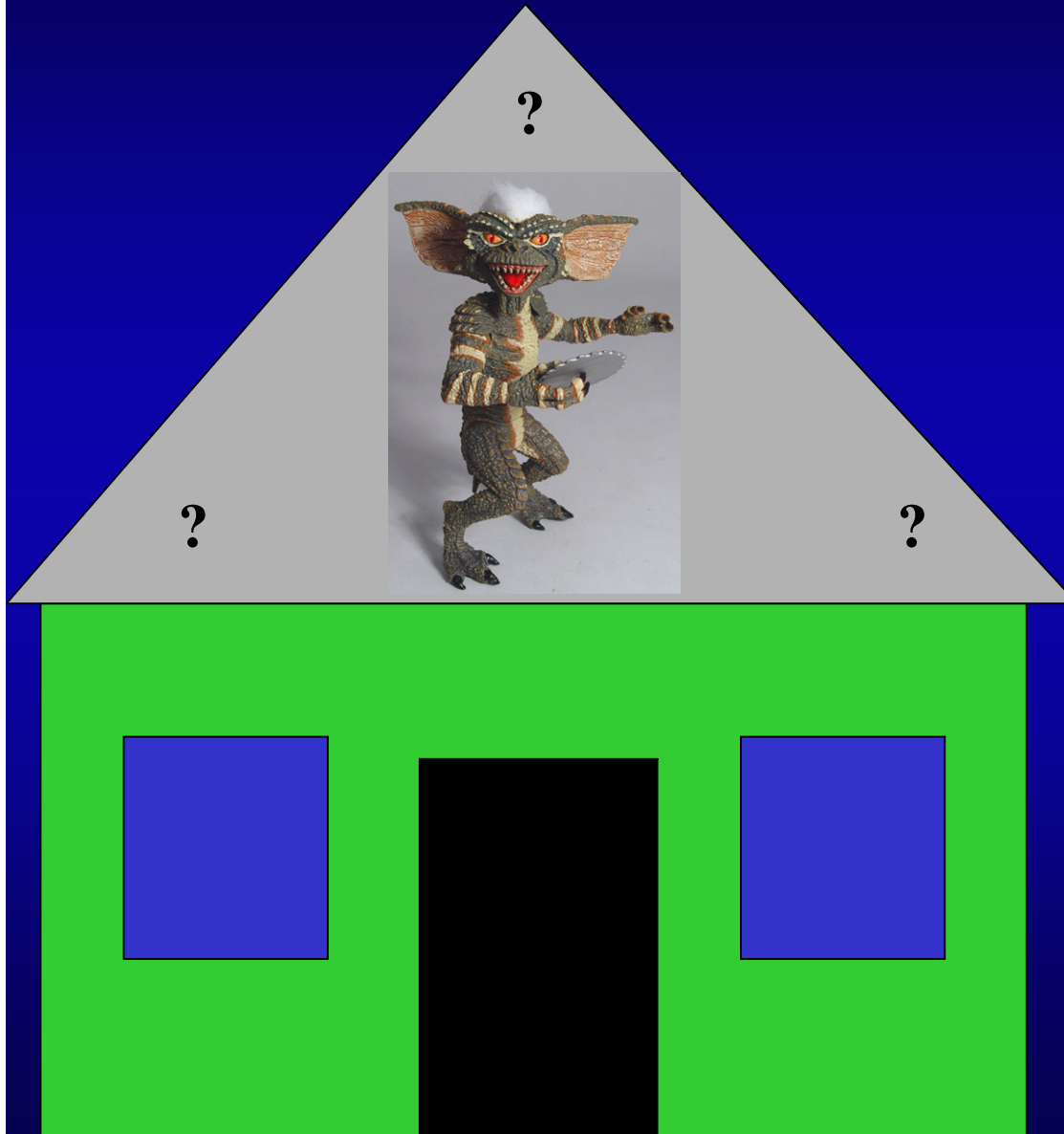
represents the degree to which we believe a given **model** accurately describes the situation given the available **data** and all of our prior information **I**

Prior  
Probability

describes the degree to which we believe the model accurately describes reality based on all of our prior information.

Normalizing  
constant

# Elliot Sober's Gremlins



Observation: Loud noise  
in the attic

Hypothesis: *gremlins in the  
attic playing bowling*

Likelihood =

$P(\text{noise} | \text{gremlins in the attic})$   
very high

Posterior Probability =

$P(\text{gremlins in the attic} | \text{noise})$   
very low

# ML Mapping

(Strimmer and von Haeseler, 1997)

**Data:** Alignment of four sequences

**Hypotheses:** All possible unrooted tree topologies

$T_1, T_2, T_3$

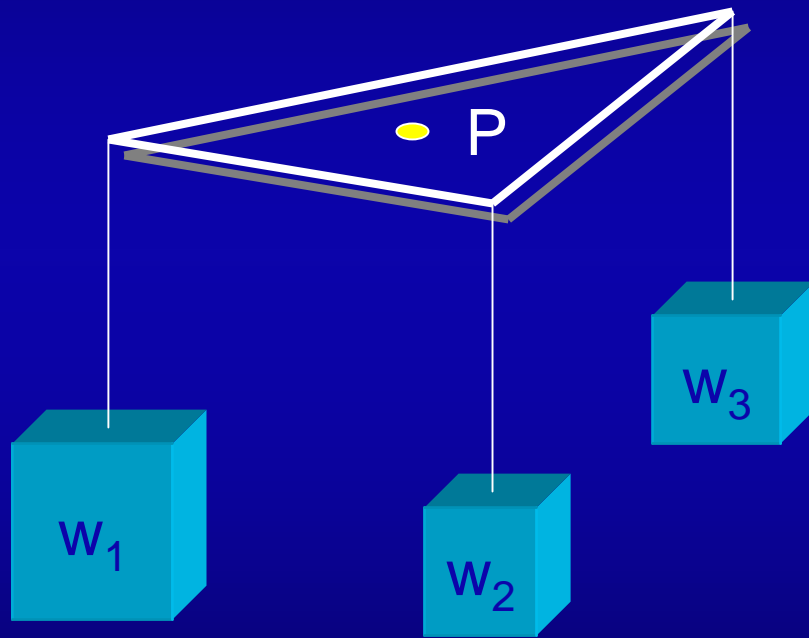
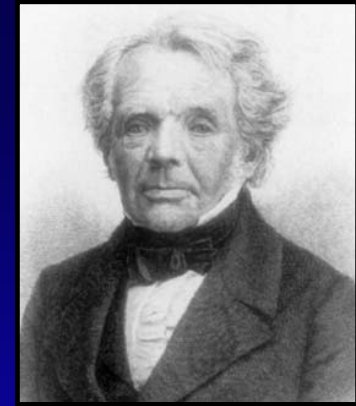
**Prior:** Equal Probabilities

For each set of 4 sequences:

- Calculate maximum-likelihood  $L_i$  for each tree  $T_i$
- Calculate posterior probabilities  $p_i$  for each tree  $T_i$
- Plot the point  $(p_1, p_2, p_3)$  into equilateral triangle

# Barycentric Coordinates

(August Ferdinand Möbius, 1827)



**P : barycenter=center of gravity**

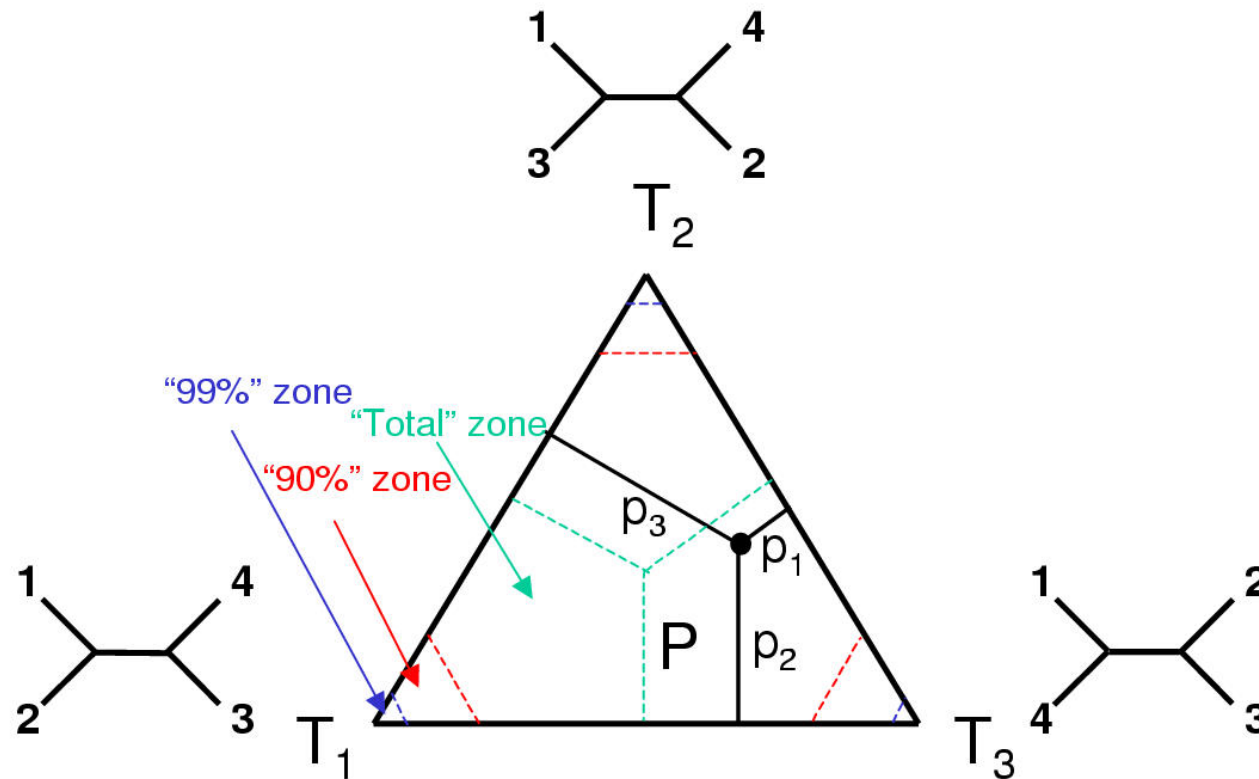
For any point P inside the triangle, there exist masses  $w_1, w_2, w_3$  such that if placed at the corresponding vertices of the triangle, their center of gravity will coincide with point P.

Barycentric coordinates are defined uniquely for every point inside the triangle (given that  $w_1+w_2+w_3=1$ ).



# ML Mapping

(Fig. modified from Strimmer)



$p_1$ ,  $p_2$  and  $p_3$  are barycentric coordinates of point  $P$

# Data Flow

Download four genomes (genome quartet) [a.a.sequences]

“BLAST” every genome against every other genome

Select top hit of every BLAST search

Detect quartets of orthologs

Align quartets of orthologues using ClustalW

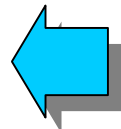
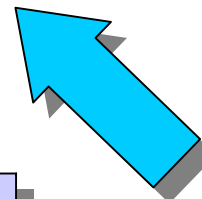
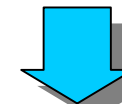
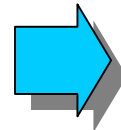
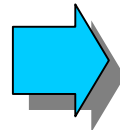
Calculate maximum-likelihood values and posterior probabilities for all three tree topologies

Convert probabilities (barycentric coordinates) into Cartesian coordinates

Plot all points onto equilateral triangle

Extract datasets with strong preference for a particular topology ( $p > 0.99$ )

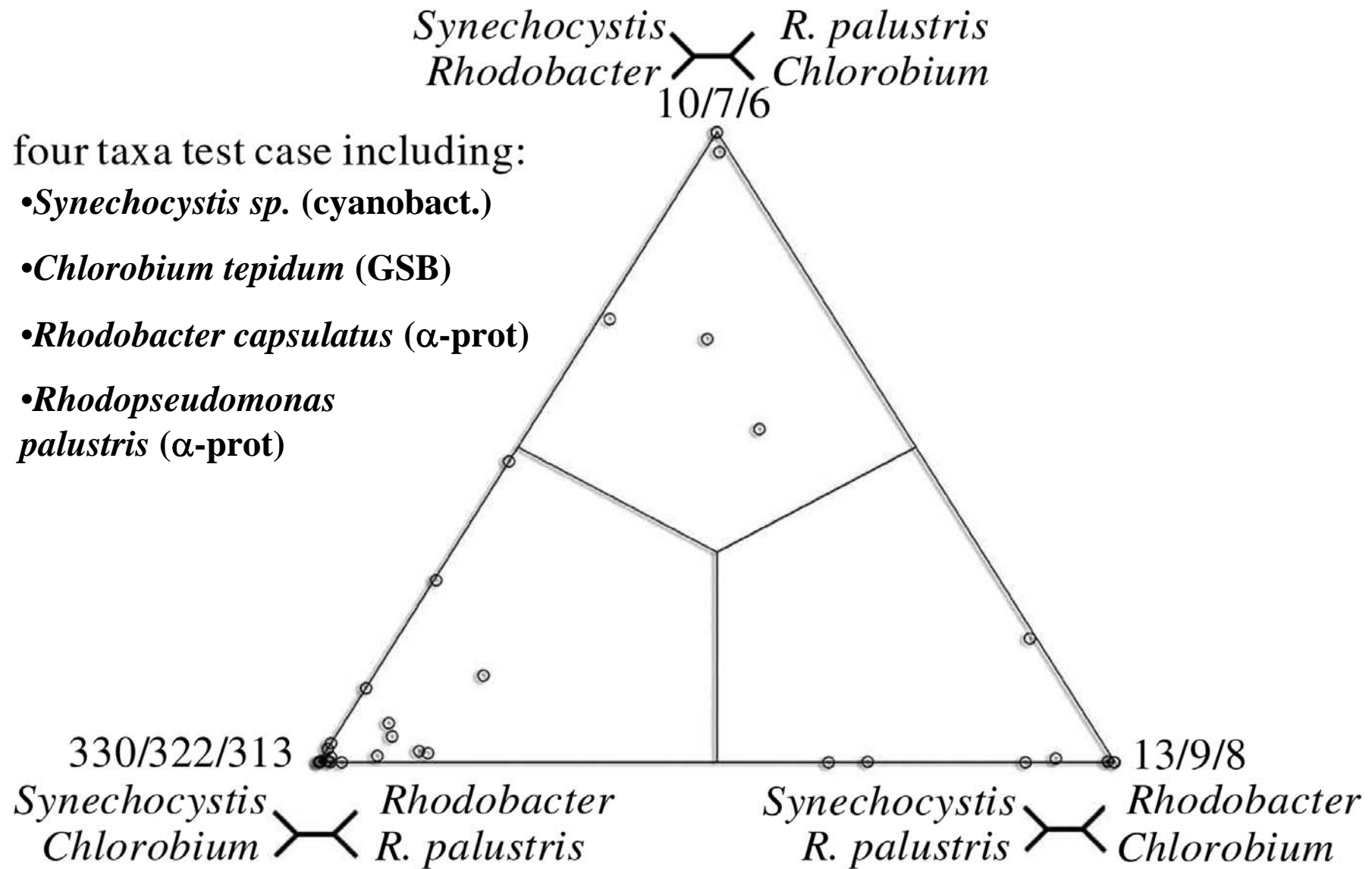
Detect Functional Category (according to COG database)



# TEST CASE

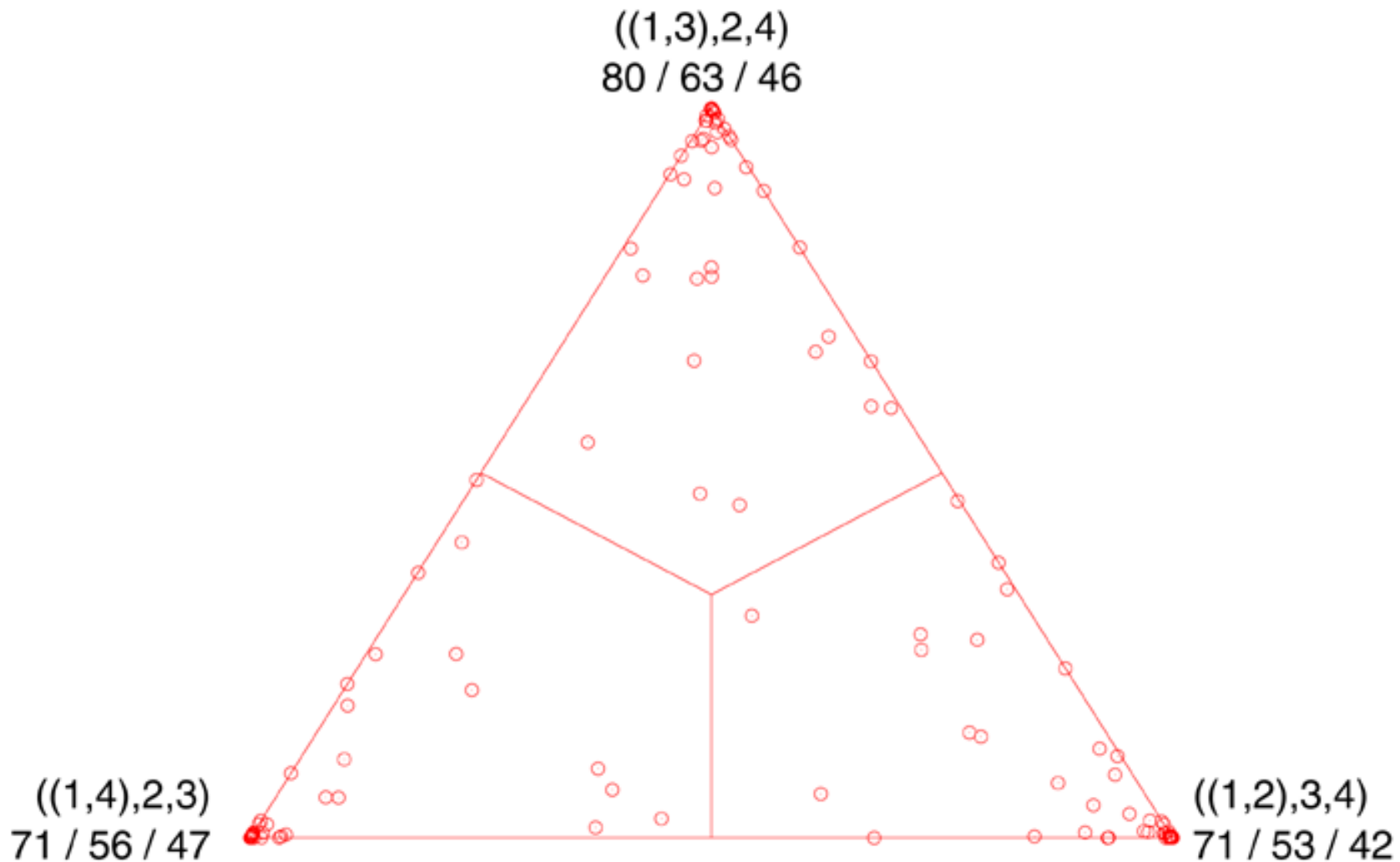
four taxa test case including:

- *Synechocystis* sp. (cyanobact.)
- *Chlorobium tepidum* (GSB)
- *Rhodobacter capsulatus* ( $\alpha$ -prot)
- *Rhodospseudomonas palustris* ( $\alpha$ -prot)



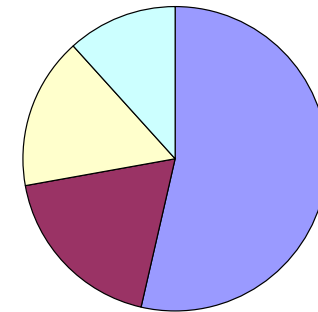
# Inter-phylum relationships (bacteria) - there is no obvious core

#8: E.coli (1), D.radiodurans (2), B. subtilis (3), T.pallidum (4)

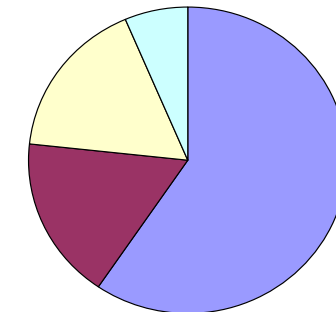


		#8		
<b>Functional Categories of COGs:</b>		<b>1</b>	<b>2</b>	<b>3</b>
<b>Information storage and processing</b> <span style="color:blue">■</span>		23	28	25
<b>J</b>	Translation, ribosomal structure and biogenesis	15	22	15
<b>K</b>	Transcription	0	0	4
<b>L</b>	DNA replication, recombination and repair	8	6	6
<b>Cellular processes</b> <span style="color:maroon">■</span>		8	8	11
<b>D</b>	Cell division and chromosome partitioning	0	2	0
<b>O</b>	Posttranslational modification, protein turnover, chaperones	4	2	4
<b>M</b>	Cell envelope biogenesis, outer membrane	3	3	1
<b>N</b>	Cell motility and secretion	1	1	5
<b>P</b>	Inorganic ion transport and metabolism	0	0	1
<b>T</b>	Signal transduction mechanisms	0	0	0
<b>Metabolism</b> <span style="color:yellow">■</span>		7	8	7
<b>C</b>	Energy production and conversion	1	1	0
<b>G</b>	Carbohydrate transport and metabolism	2	2	3
<b>E</b>	Amino acid transport and metabolism	2	1	1
<b>F</b>	Nucleotide transport and metabolism	0	2	1
<b>H</b>	Coenzyme metabolism	2	1	2
<b>I</b>	Lipid metabolism	0	1	0
<b>Poorly characterized</b> <span style="color:cyan">■</span>		5	3	6
<b>R</b>	General function prediction only	5	3	3
<b>S</b>	Function unknown	0	0	3

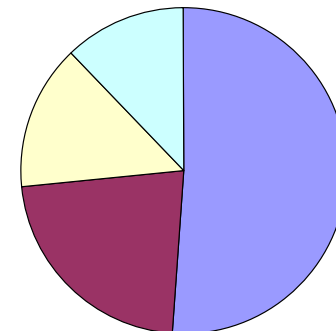
Tree #1



Tree #2



Tree #3





# Alternative Approaches to Estimate Posterior Probabilities

## Bayesian Posterior Probability Mapping with MrBayes (Huelsenbeck and Ronquist, 2001)

### Problem:

Strimmer's formula 
$$p_i = \frac{L_i}{L_1 + L_2 + L_3}$$
 only considers 3 trees  
(those that maximize the likelihood for the three topologies)

### Solution:

Exploration of the tree space by sampling trees using a biased random walk  
(Implemented in MrBayes program)

Trees with higher likelihoods will be sampled more often

$$p_i \approx \frac{N_i}{N_{\text{total}}}$$

,where  $N_i$  - number of sampled trees of topology  $i$ ,  $i=1,2,3$

$N_{\text{total}}$  - total number of sampled trees (has to be large)

# Illustration of a biased random walk

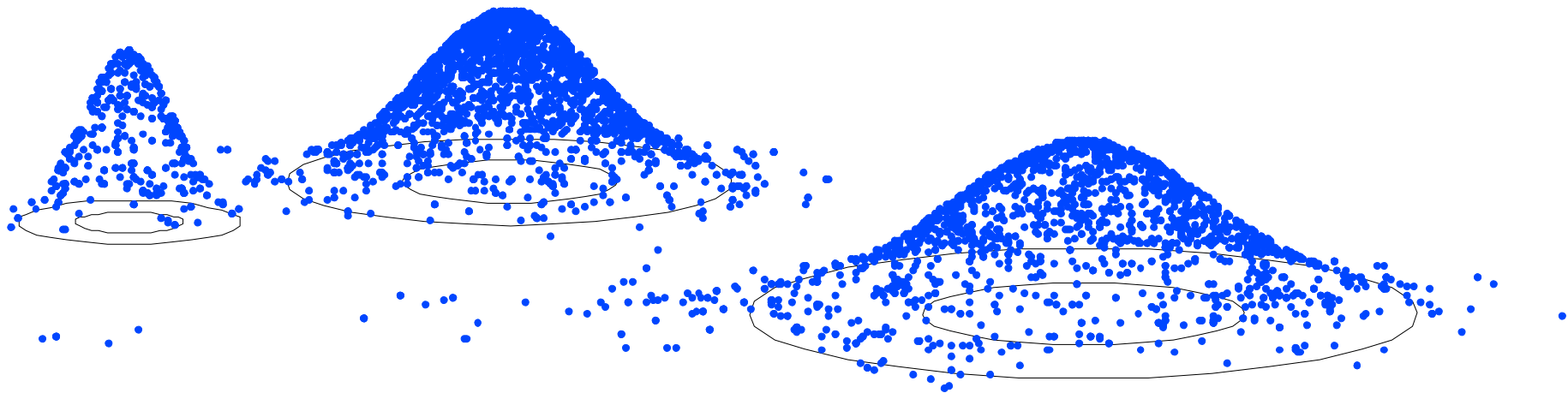
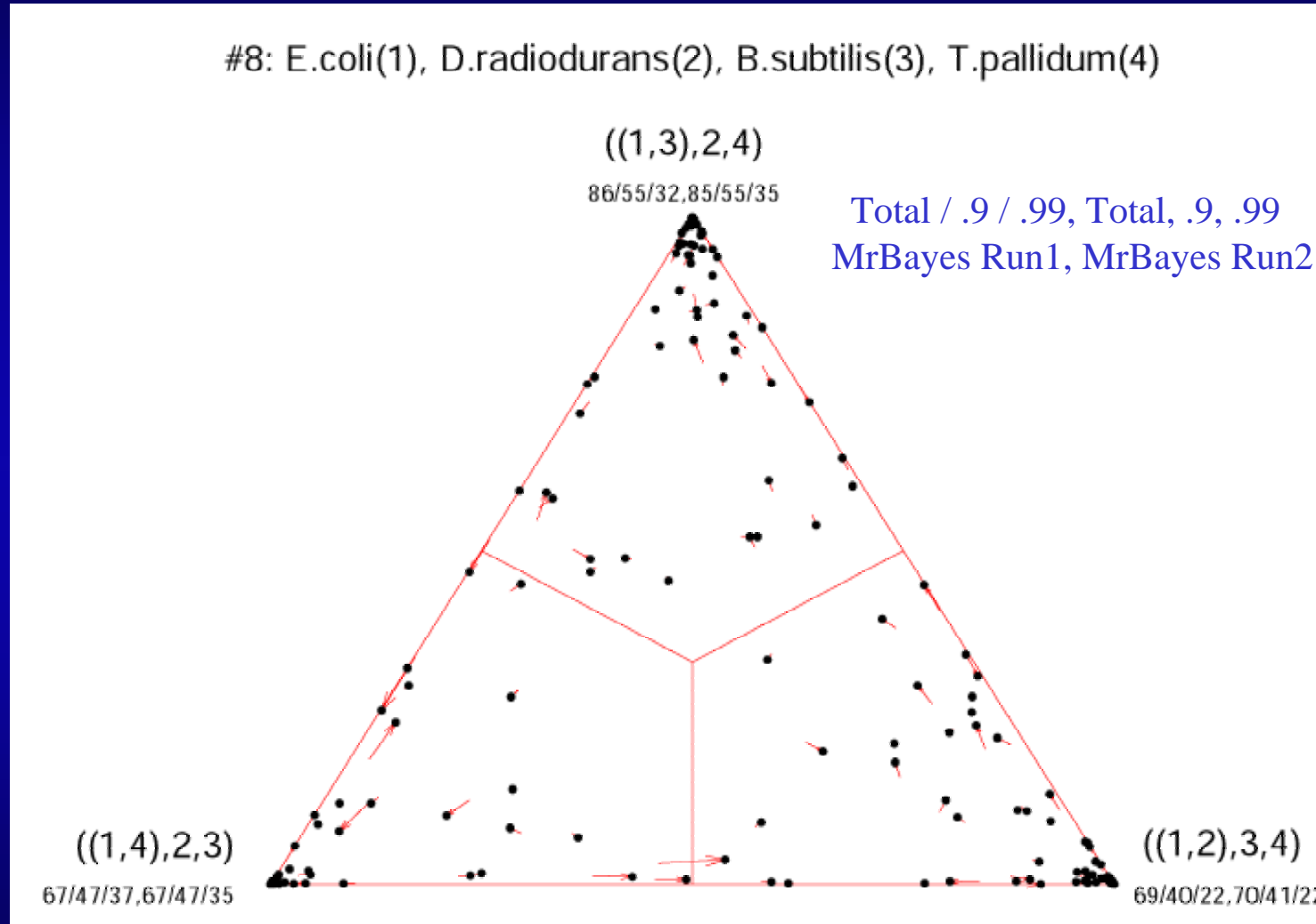


Figure generated using MCRobot program (Paul Lewis, 2001)

# Inter-phylum relationships (bacteria) - there is no obvious core

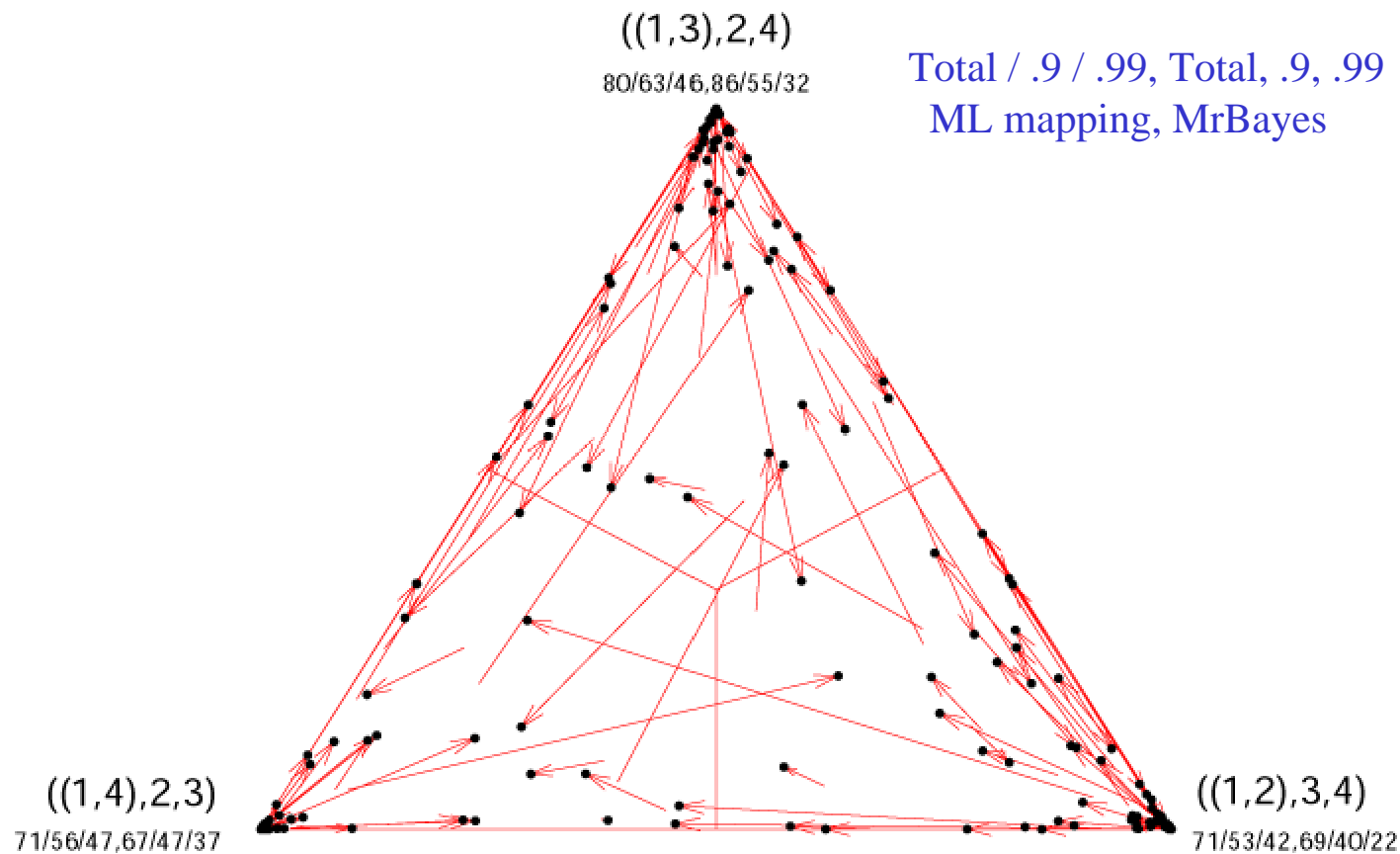


P-vector with MrBayes Run#1: Start of arrow

P-vector with MrBayes Run#2: Black dot at tip of arrow

# Comparing ML-mapping to Bayesian posterior probabilities

#8: E.coli(1), D.radiodurans(2), B.subtilis(3), T.pallidum(4)



P-vector with ML-mapping: Start of arrow

P-vector with MrBayes: Black dot at tip of arrow

# Alternative Approaches to Estimate Posterior Probabilities (2)

Bootstrap Support Values Mapping:

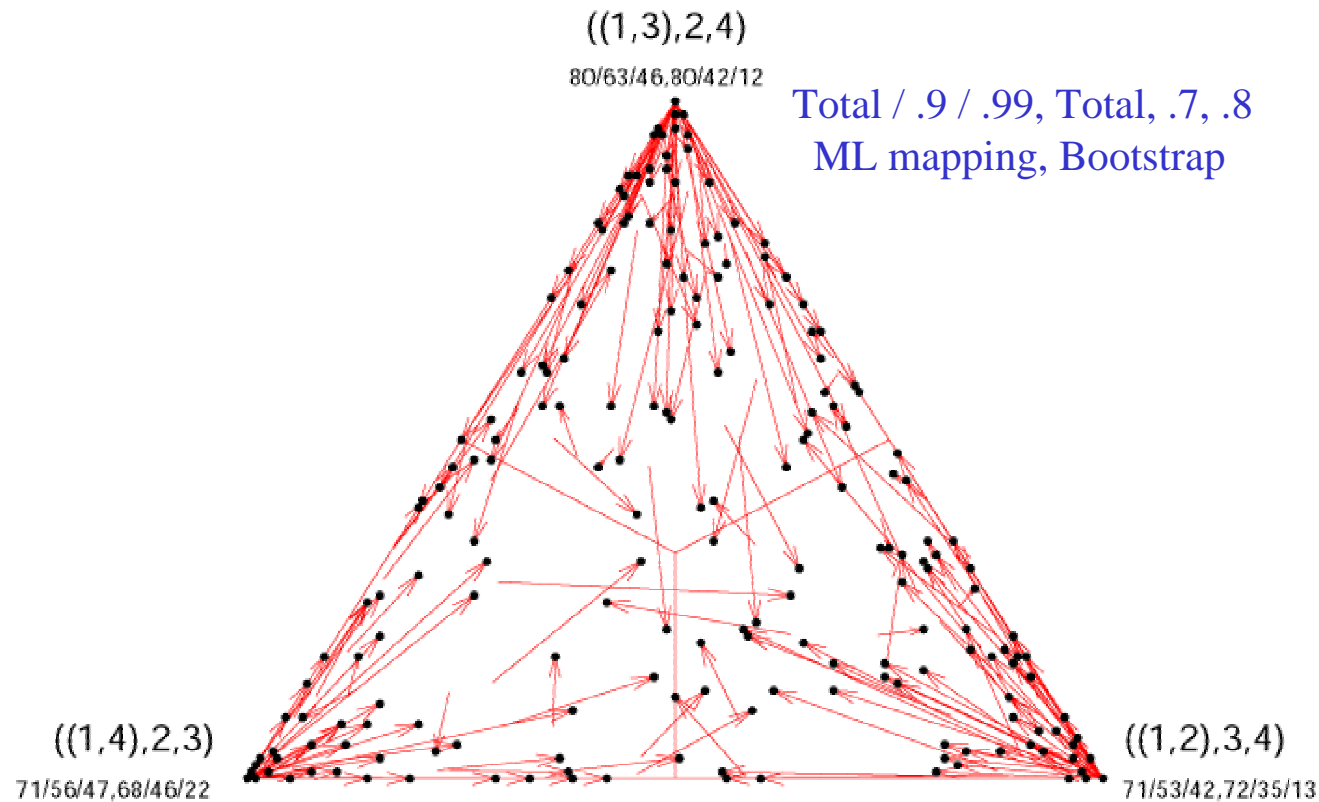
For each Quartet of Orthologous Proteins:

- 1) Create 100 bootstrapped samples
- 2) Evaluate three tree topologies for each of 100 samples
- 3) Construct bootstrap support values vector, i.e., percent of bootstrapped samples that have the highest likelihood value for each tree topology.



# Comparing ML-Mapping to Bootstrap Support Values

#8: E.coli(1), D.radiodurans(2), B.subtilis(3), T.pallidum(4)



P-vector with ML-mapping: Start of arrow

P-vector with Bootstrap: Black dot at tip of arrow

# Comparing Support Measures:

99%  $\approx$  90%  $\approx$  70%

posterior probability  
calculated according  
to **ML mapping**

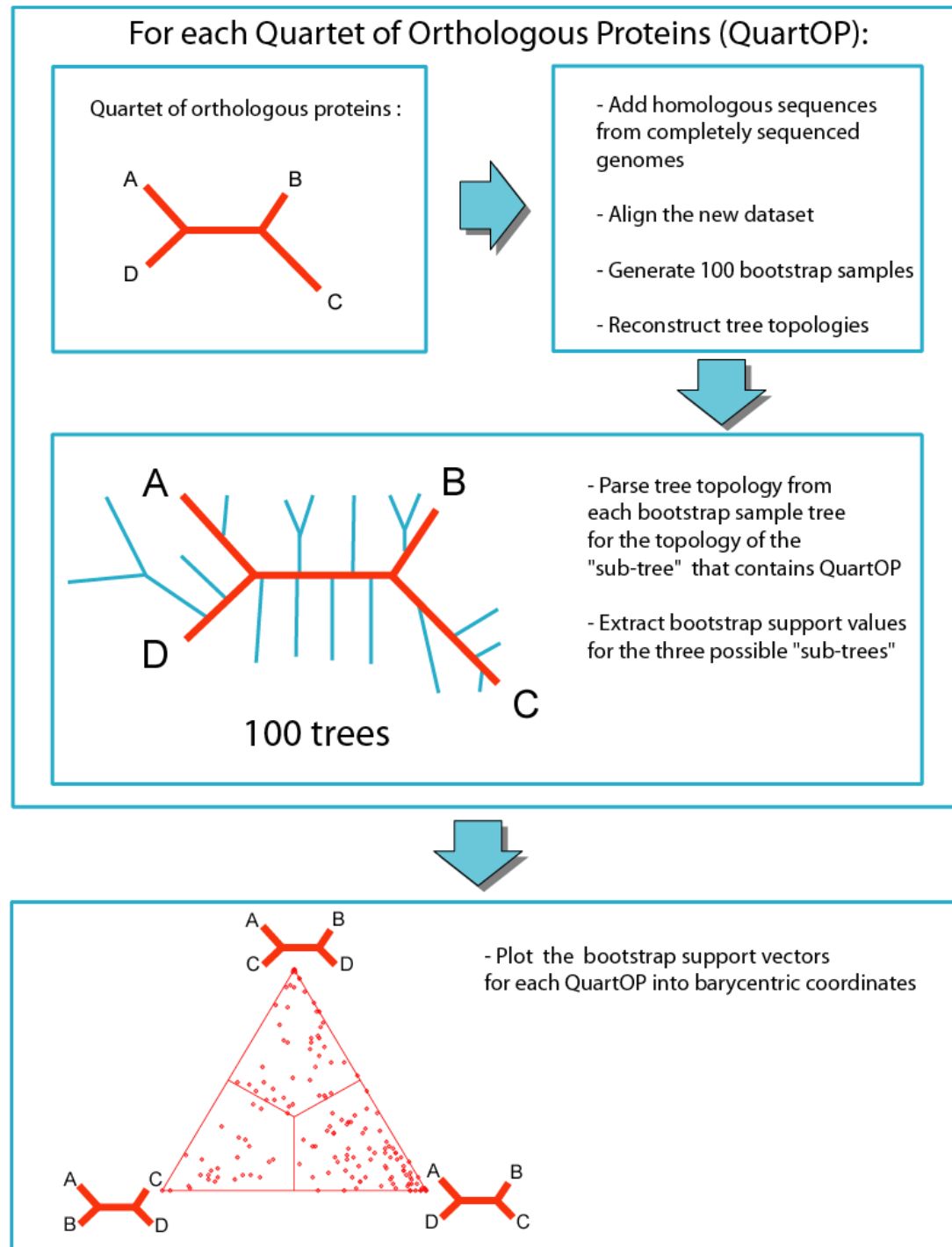
posterior probability  
estimated using  
**MCMC** (MrBayes)

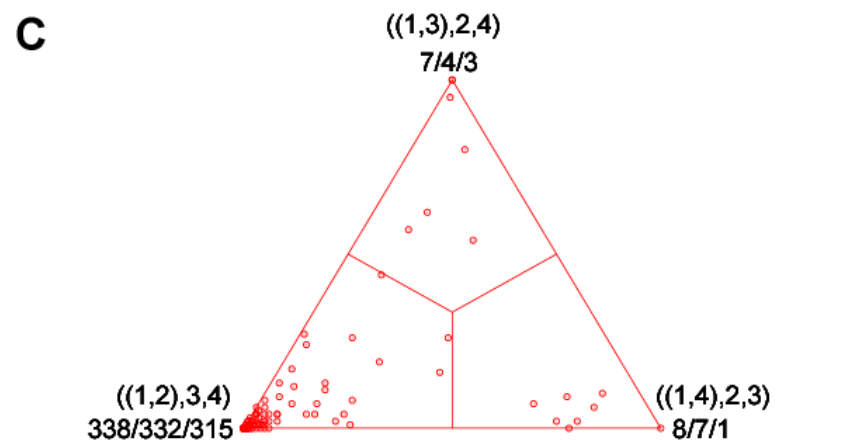
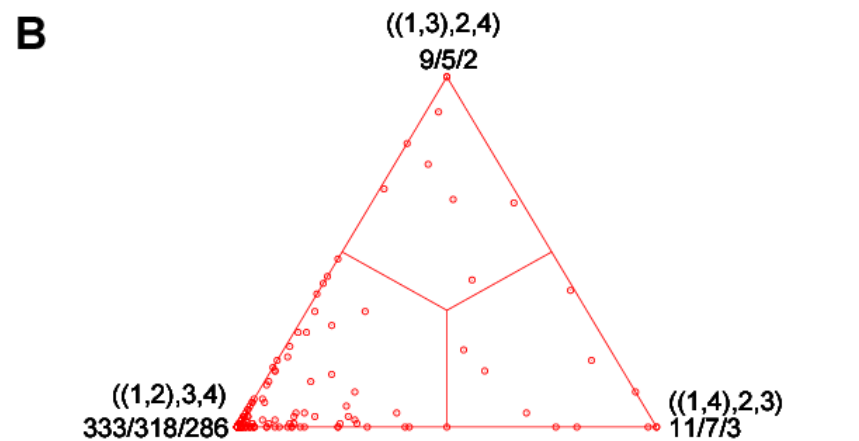
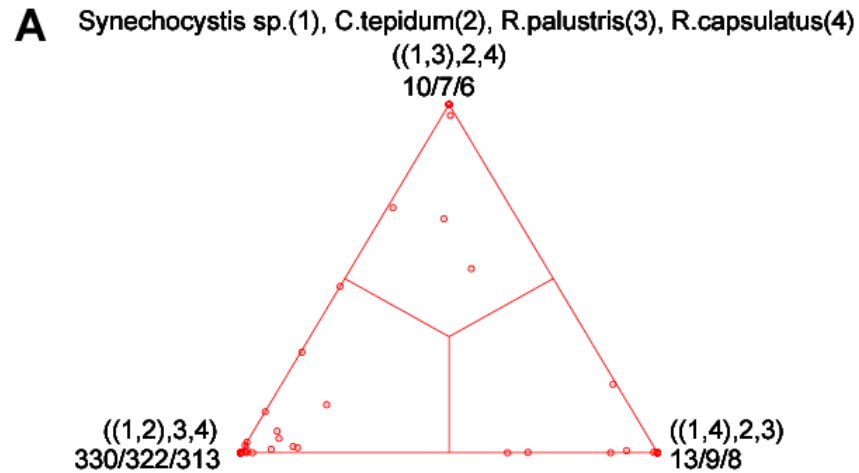
**bootstrap**  
support

# Increasing Reliability

Phylogenetic reconstruction becomes more reliable when more sequences are included.

**DATA FLOW**  
analyses of  
extended  
datasets





## COMPARISON OF DIFFERENT SUPPORT MEASURES

**A:** mapping of posterior probabilities according to Strimmer and von Haeseler

**B:** mapping of bootstrap support values

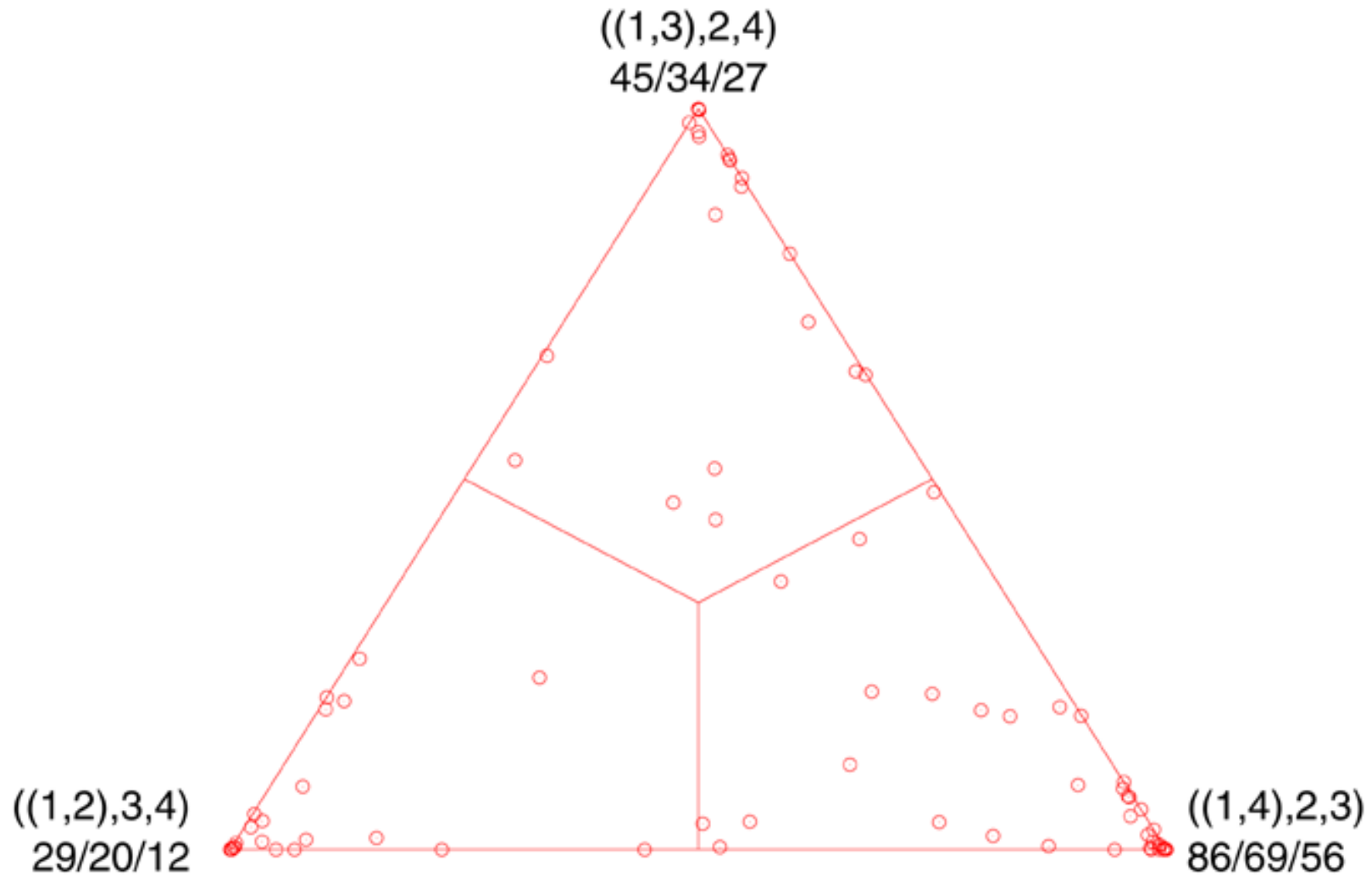
**C:** mapping of bootstrap support values from extended datasets

# Inter-Domain Genome Comparisons

- ✓ *Synechocystis* sp. – cyanobacterium
- ✓ *Thermotoga maritima* – thermophilic bacterium
- ✓ *Aquifex aeolicus* – thermophilic bacterium
- ✓ *Halobacterium* sp. – salt-loving euryarchaeon

# ML Map

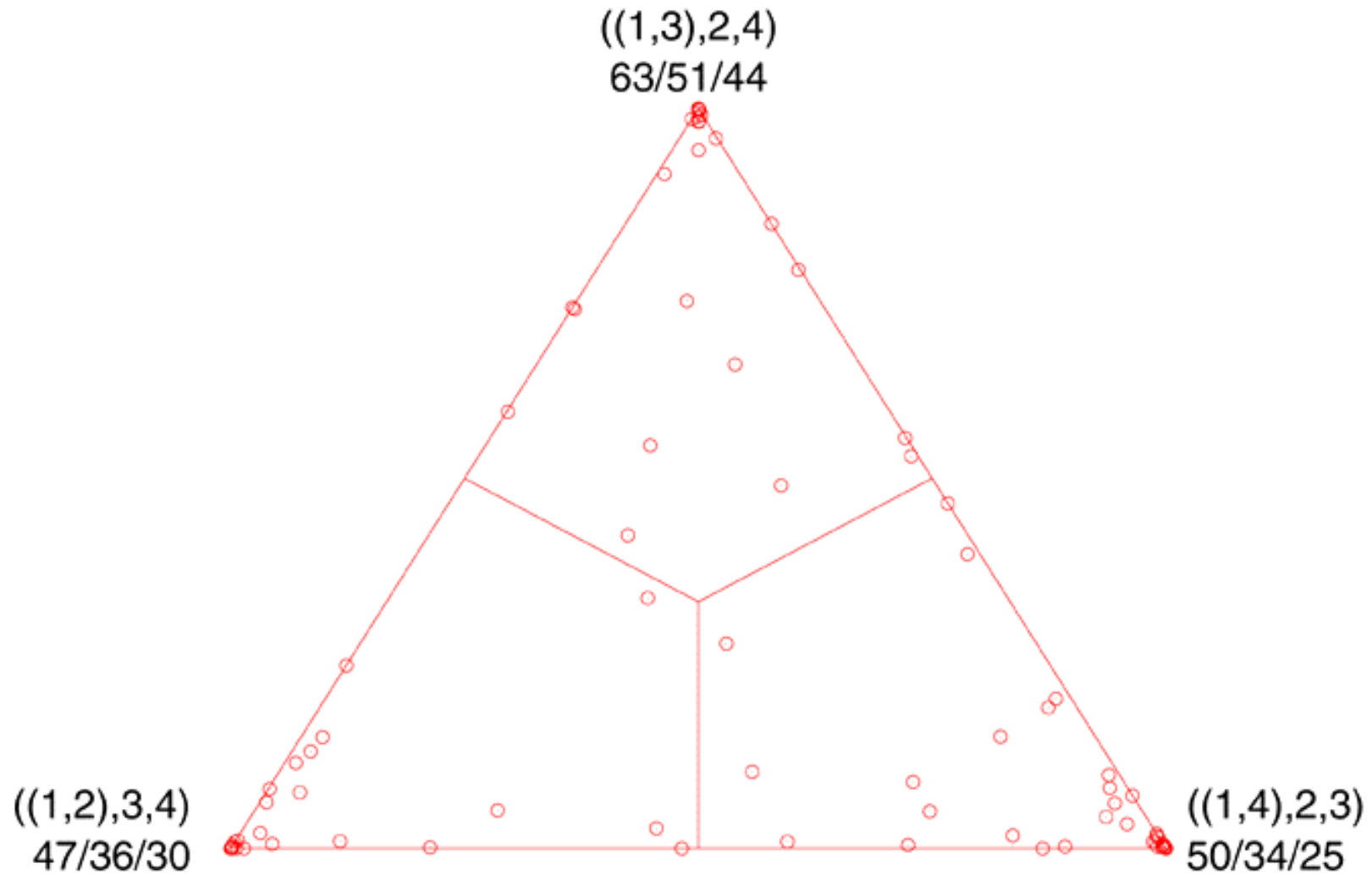
#11: *Synechocystis* sp.(1), *T.maritima*(2), *A.aeolicus*(3), *Halobacterium* sp.(4)





# ML Map

#13: *Synechocystis* sp.(1), *T.maritima*(2), *A.aeolicus*(3), *A.fulgidus*(4)

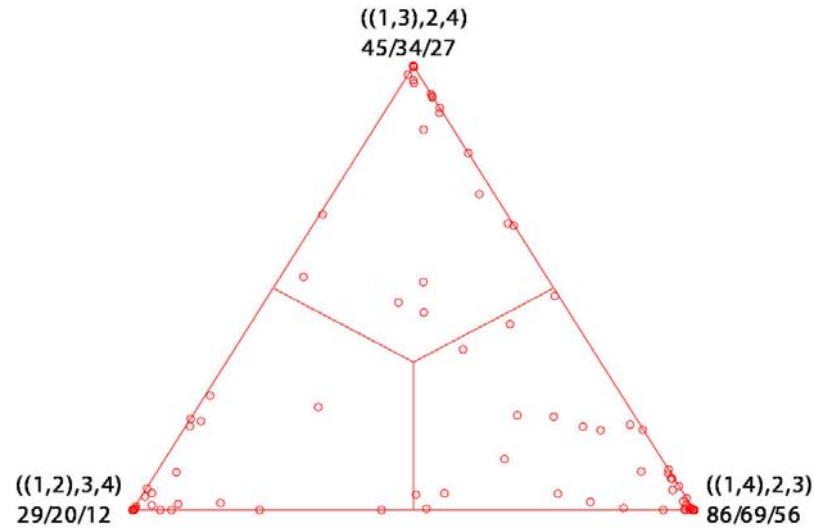


ml-mapping

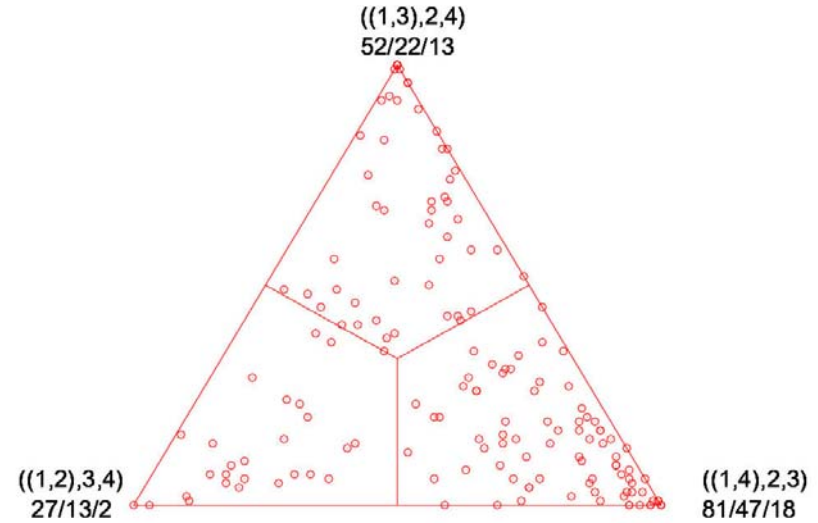
versus

bootstrap values from  
extended datasets

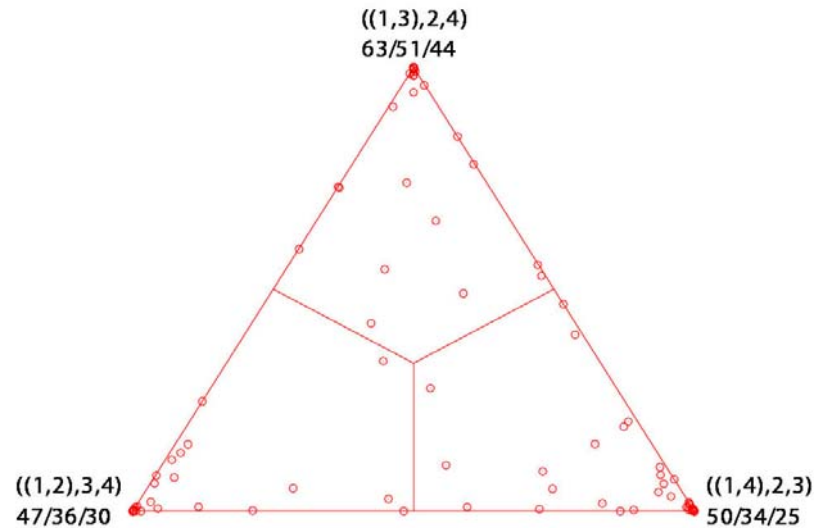
A Synechocystis sp.(1), T.maritima(2), A.aeolicus(3), Halobacterium sp.(4)



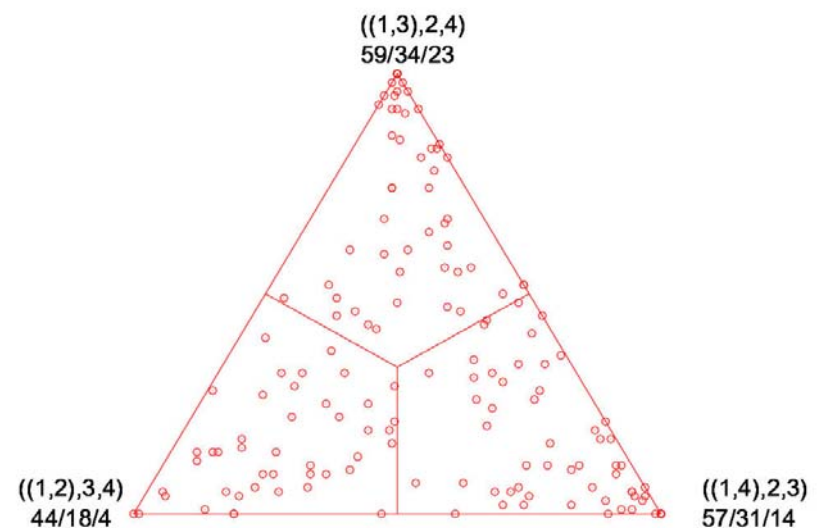
A Synechocystis sp.(1), T.maritima(2), A.aeolicus(3), Halobacterium sp.(4)



B Synechocystis sp.(1), T.maritima(2), A.aeolicus(3), A.fulgidus(4)



B Synechocystis sp.(1), T.maritima(2), A.aeolicus(3), A.fulgidus(4)



Proteins in the *information storage and processing category* that group orthologs from *Halobacterium* with *Synechocystis* and *Thermotoga* with *Aquifex* (Topology #3 – putative identification)

### Nucleotide modifying Enzymes

- tRNA-pseudouridine synthase
- dimethyladenosine transferase

### Enzymes involved in DNA repair and recombination

- DNA mismatch repair protein
- excision nuclease A,B,C chains (involved in DNA repair)
- Endonuclease V (involved in DNA repair)

### Enzymes involved in translation

- putative translation factor SUA5
- initiation factor IF2
- translation initiation factor eIF-2B subunit alpha
- Glu-tRNA amidotransferase subunits A,B
- ribosomal proteins L1,L11,L3,S4
- amino acyl tRNA synthetases for serine, valine, methionine, cysteine, proline, phenylalanine ( $\alpha$  SU)

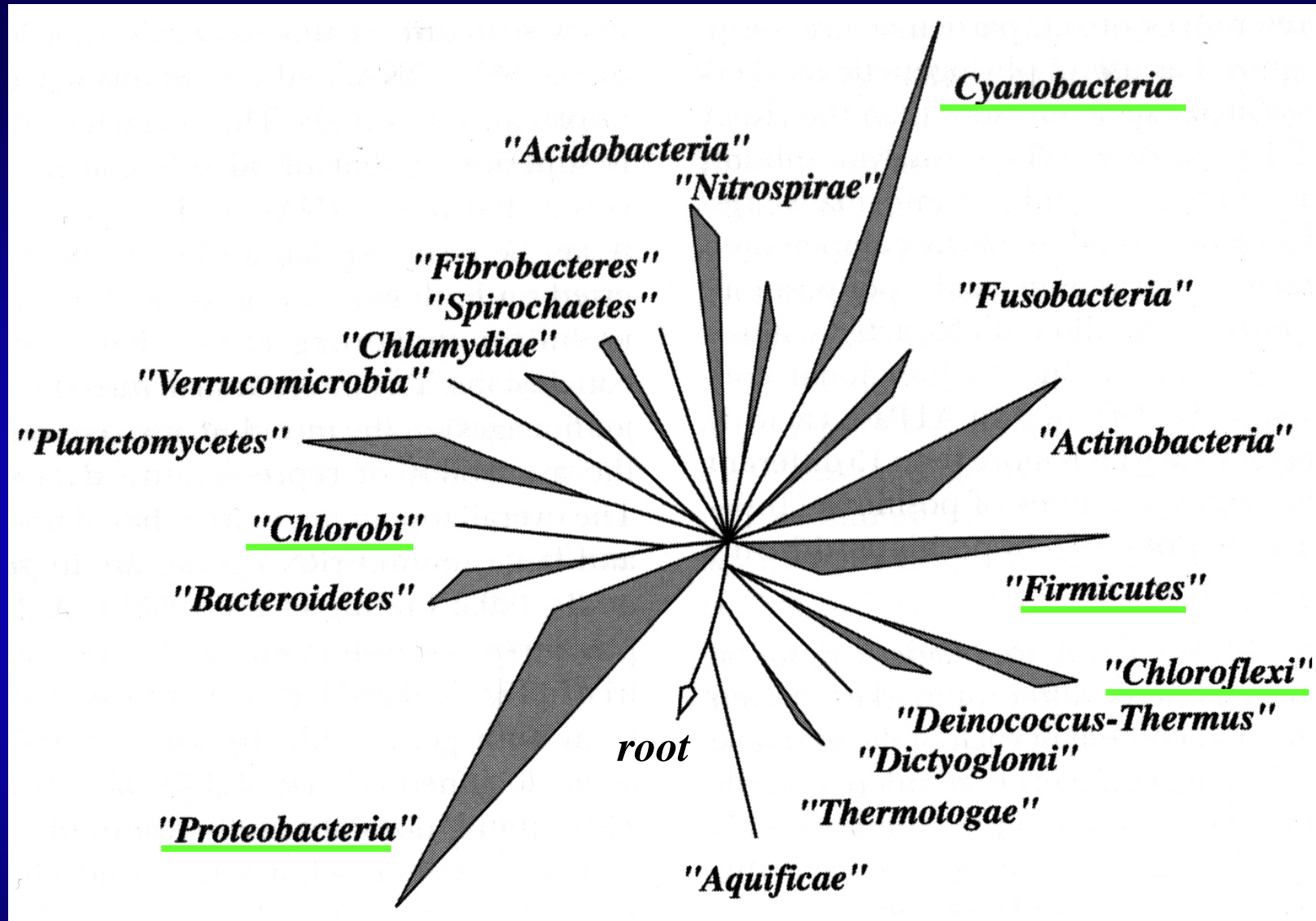
### Other

- DNA gyrase subunits [A,B]
- DNA helicase

## NUMBER OF GENES PER CONFIDENCE LEVEL FOR DIFFERENT TYPES OF MAPPINGS

Genome Quartet	99% posterior probability	90% bootstrap support from non-extended datasets	90% bootstrap support from extended datasets
Interdomain quartet consisting of <i>Synechocystis</i> sp., <i>Halobacterium</i> sp., <i>Aquifex aeolicus</i> and <i>Thermotoga maritima</i> .	95	42	33
Interdomain quartet consisting of <i>Synechocystis</i> sp., <i>Archaeoglobus fulgidus</i> , <i>Aquifex aeolicus</i> and <i>Thermotoga maritima</i> .	99	42	41
Interphylum quartet of <i>Synechocystis</i> sp., <i>Chlorobium tepidum</i> , <i>Rhodobacter capsulatus</i> and <i>Rhodopseudomonas palustris</i> .	327	291	319

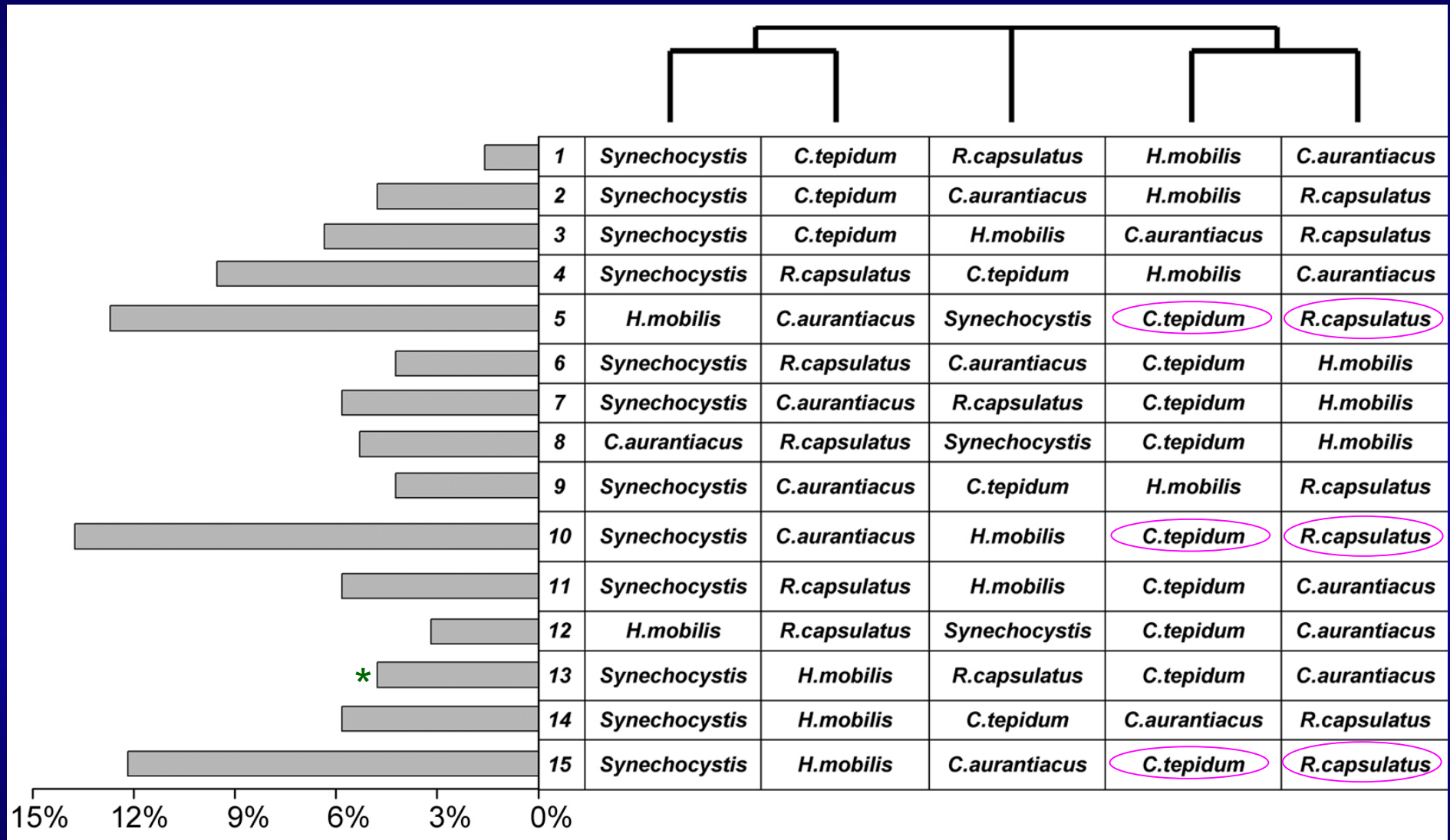
# **Extension of Mapping to Five Genomes**



## 23S rRNA tree depicting the major bacterial phyla

(from Bergey's Manual of Systematic Bacteriology, 2<sup>nd</sup> Ed.)

# Distribution of orthologs among 15 possible trees

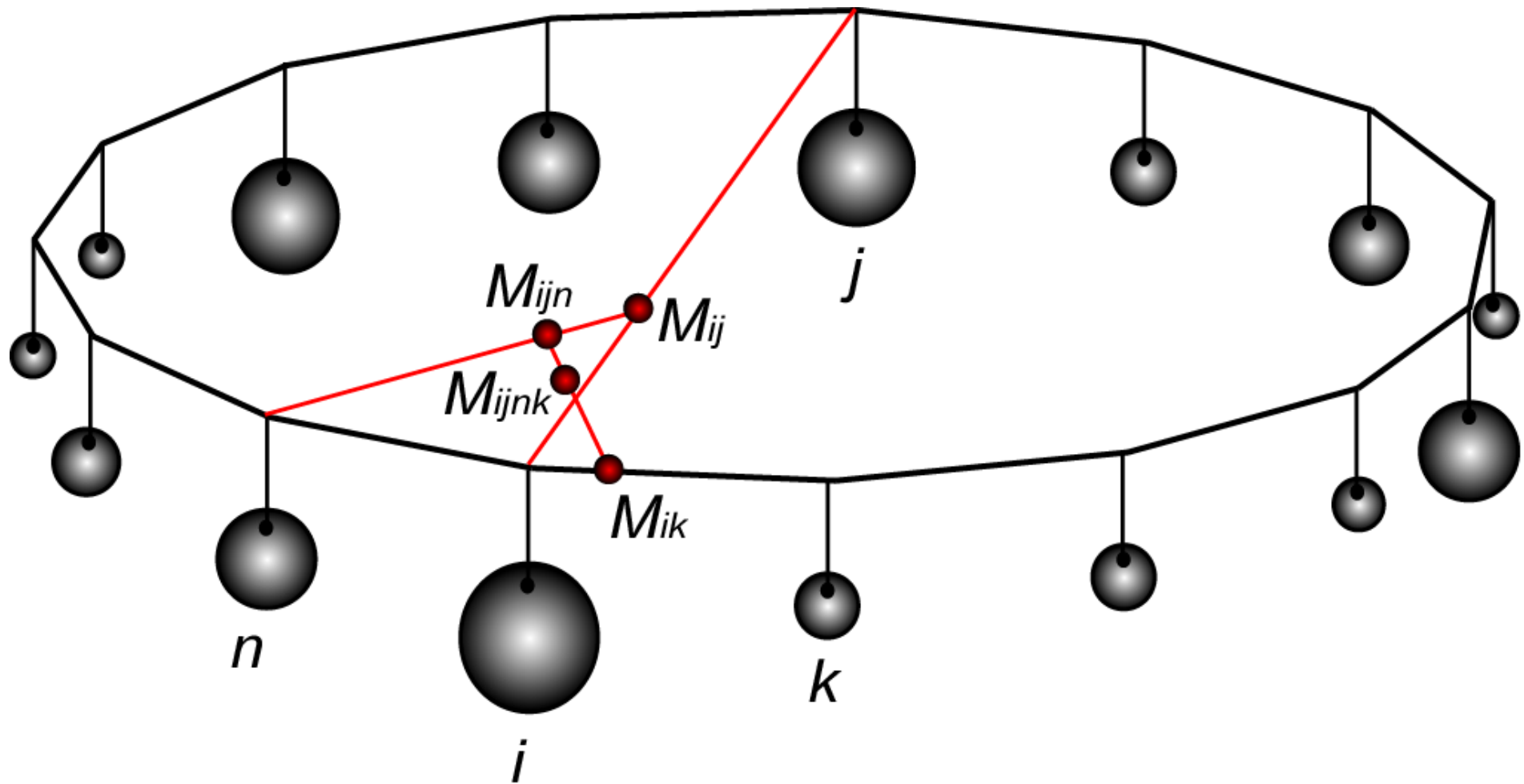


188 datasets of orthologous genes



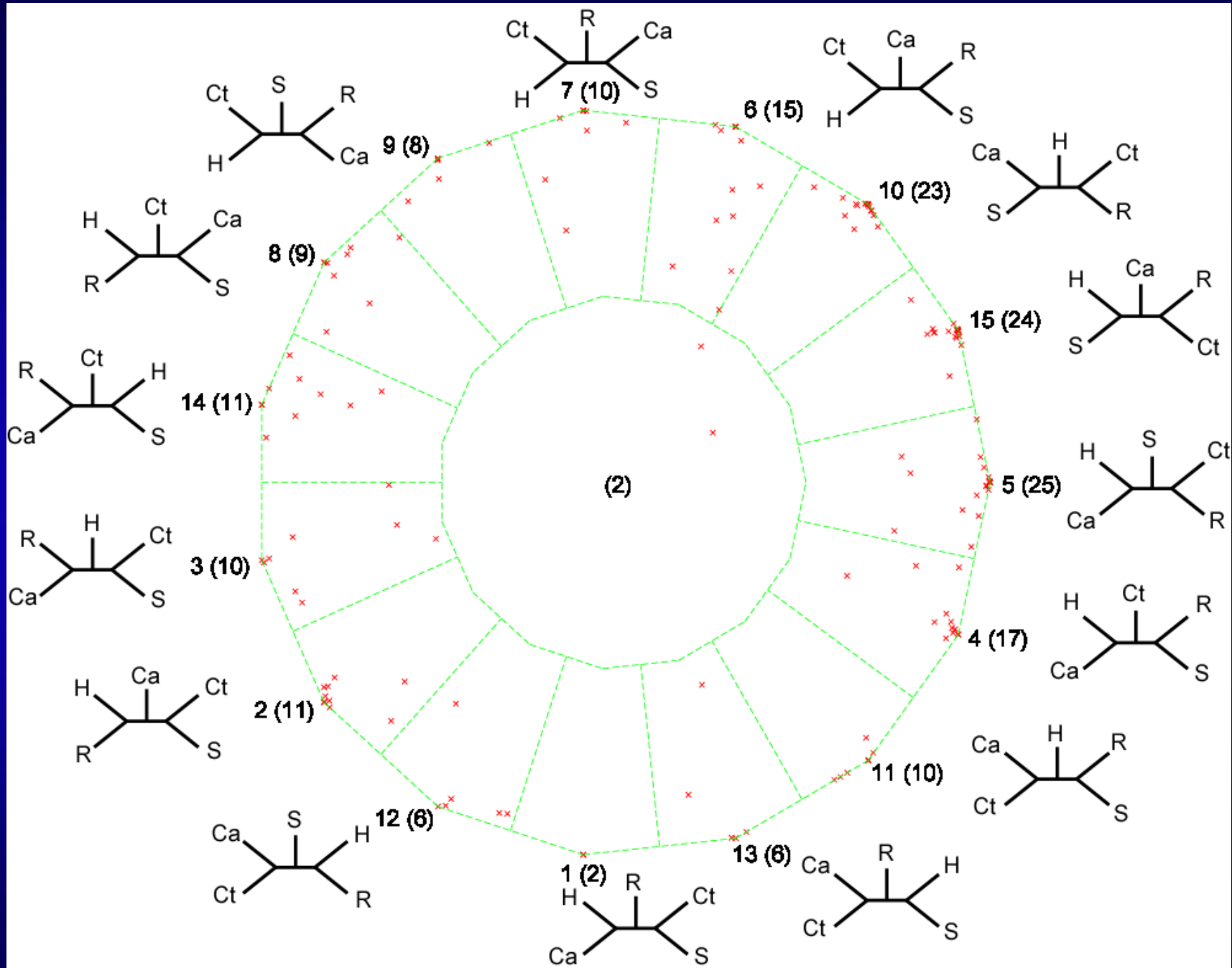
# CALCULATION OF THE CENTER OF GRAVITY OF THE DEKAPENTAGON

Illustration of the principle



# PHYLOGENETIC DEKAPENTAGON

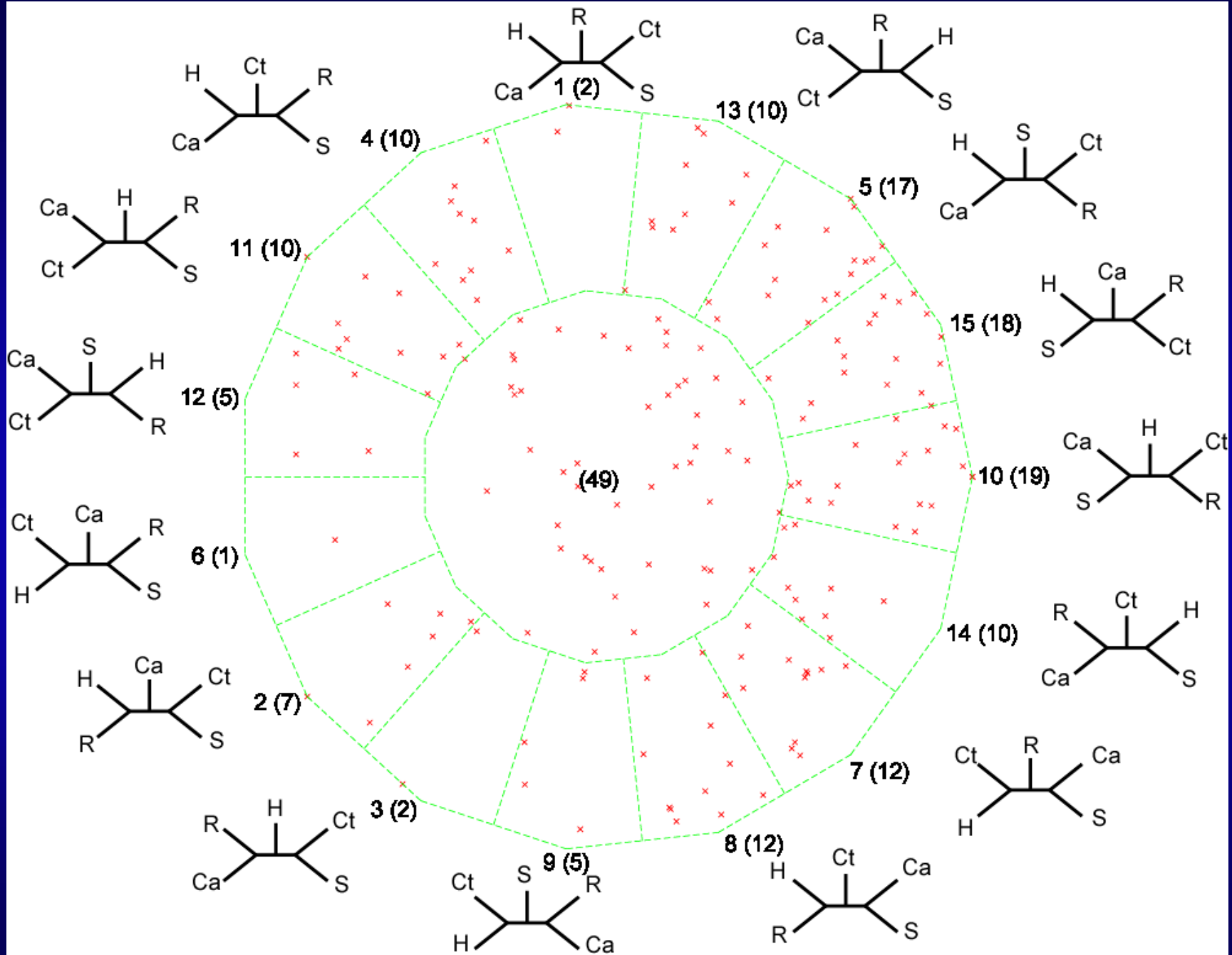
Posterior probabilities



R: *Rhodobacter capsulatus*, H: *Heliobacillus mobilis*, S: *Synechocystis* sp., Ct: *Chlorobium tepidum*, Ca: *Chloroflexus aurantiacus*

# PHYLOGENETIC DEKAPENTAGON

Bootstrap support values



# Extension of the analyses to more than five genomes

## PROBLEM:

Number of possible unrooted tree topologies is equal to  $(2n-5)!/[2^{n-3}(n-3)!]$

⇒ Polygon becomes a circle

⇒ Many topologies are not supported by data

## SOLUTION:

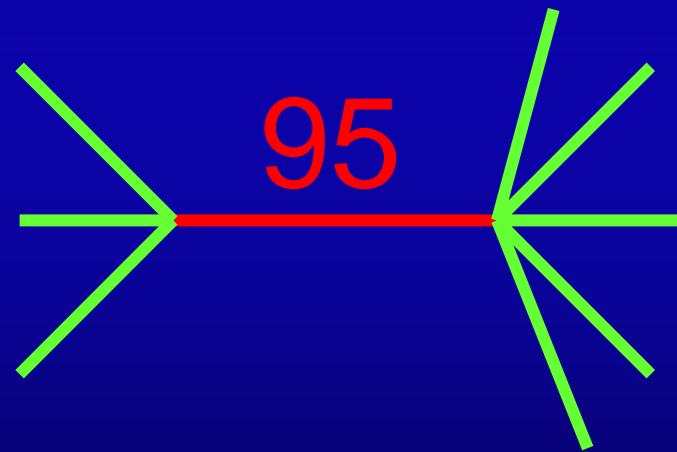
Switching from topologies to bipartitions of data

# **BIPARTITION PLOTS**

**(Modified Lento Plots)**

# BIPARTITION OF A PHYLOGENETIC TREE

**Bipartition** – a division of a phylogenetic tree into two parts that are connected by a single branch. It divides a dataset into two groups, but it does not consider the relationships within each of the two groups.



Number of bipartitions for  $N$  genomes is equal to  $2^{(N-1)} - N - 1$ .

# WHY BIPARTITIONS?

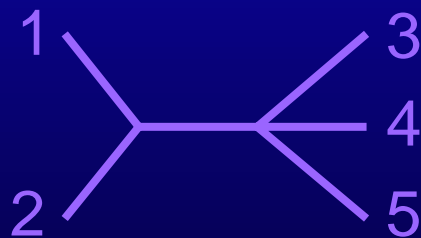
1. The number of possible bipartitions is much smaller than number of possible tree topologies, which makes it possible to evaluate all possible partitions.
2. Analyses of bipartitions allows to consider datasets that otherwise would be considered as non-informative due to lack of resolution in one or the other part of the tree.
3. Putatively horizontally transferred genes can be detected because they give rise to partitions significantly conflicting with plurality partitions.



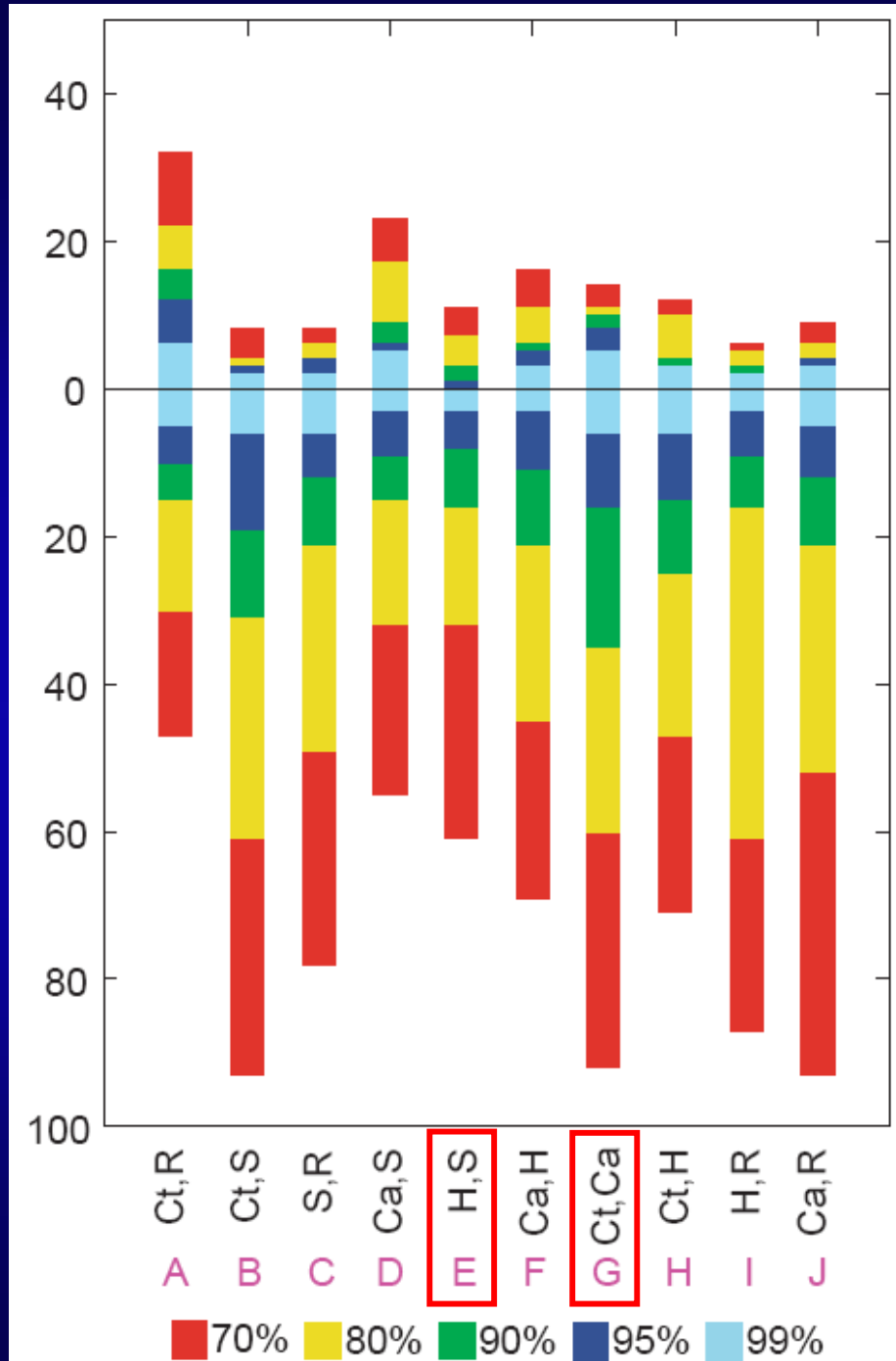
# Example of bipartition analysis for five genomes of photosynthetic bacteria

R: *Rhodobacter capsulatus*,  
 H: *Heliobacillus mobilis*,  
 S: *Synechocystis* sp.,  
 Ct: *Chlorobium tepidum*,  
 Ca: *Chloroflexus aurantiacus*

Bipartitions supported by genes from chlorophyll biosynthesis pathway

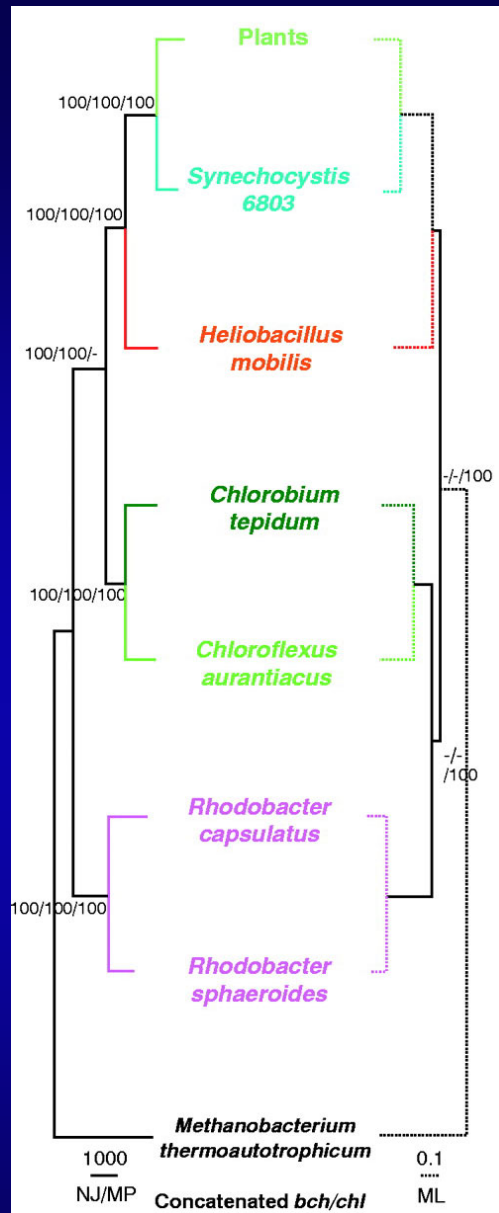


10 bipartitions



# Phylogenetic Analyses of Genes from chlorophyll biosynthesis pathway

(extended datasets)



R: *Rhodobacter capsulatus*, H: *Heliobacillus mobilis*, S: *Synechocystis sp.*, Ct: *Chlorobium tepidum*, Ca: *Chloroflexus aurantiacus*

Xiong et al. Science, 2000 289:1724-30

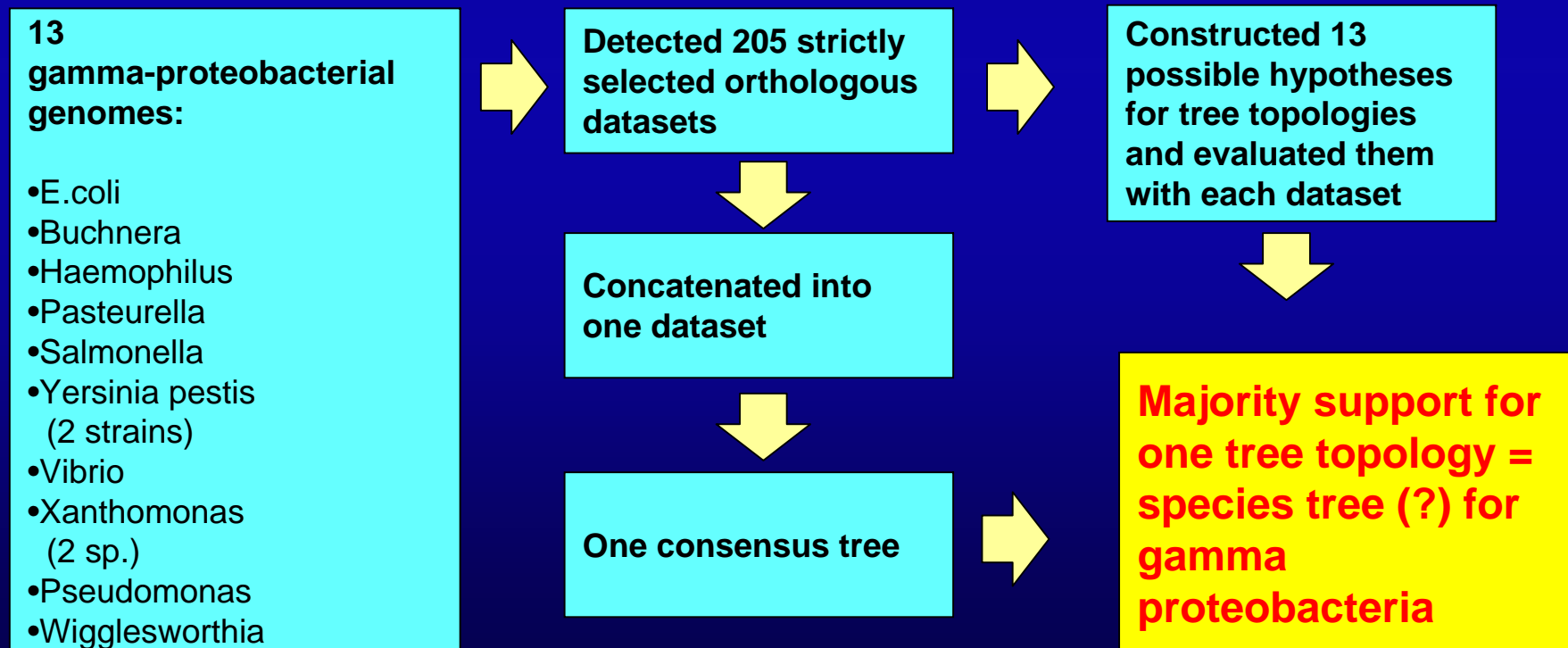
Gene	Plurality Partition		
BchB			
BchD			
BchH			
BchI			
BchL			
BchN			

TREE-PUZZLE Distances/NJ tree  
 TREE-PUZZLE Distances/FITCH tree  
 Parsimony with bootstrap  
 MrBayes (3 independent runs)

# From Gene Trees to Organismal Phylogeny in Prokaryotes: The Case of the $\gamma$ -Proteobacteria

Emmanuelle Lerat<sup>1</sup>, Vincent Daubin<sup>2</sup>, Nancy A. Moran<sup>1\*</sup>

<sup>1</sup> Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, Arizona, United States of America, <sup>2</sup> Department of Biochemistry and Molecular Biophysics, University of Arizona, Tucson, Arizona, United States of America

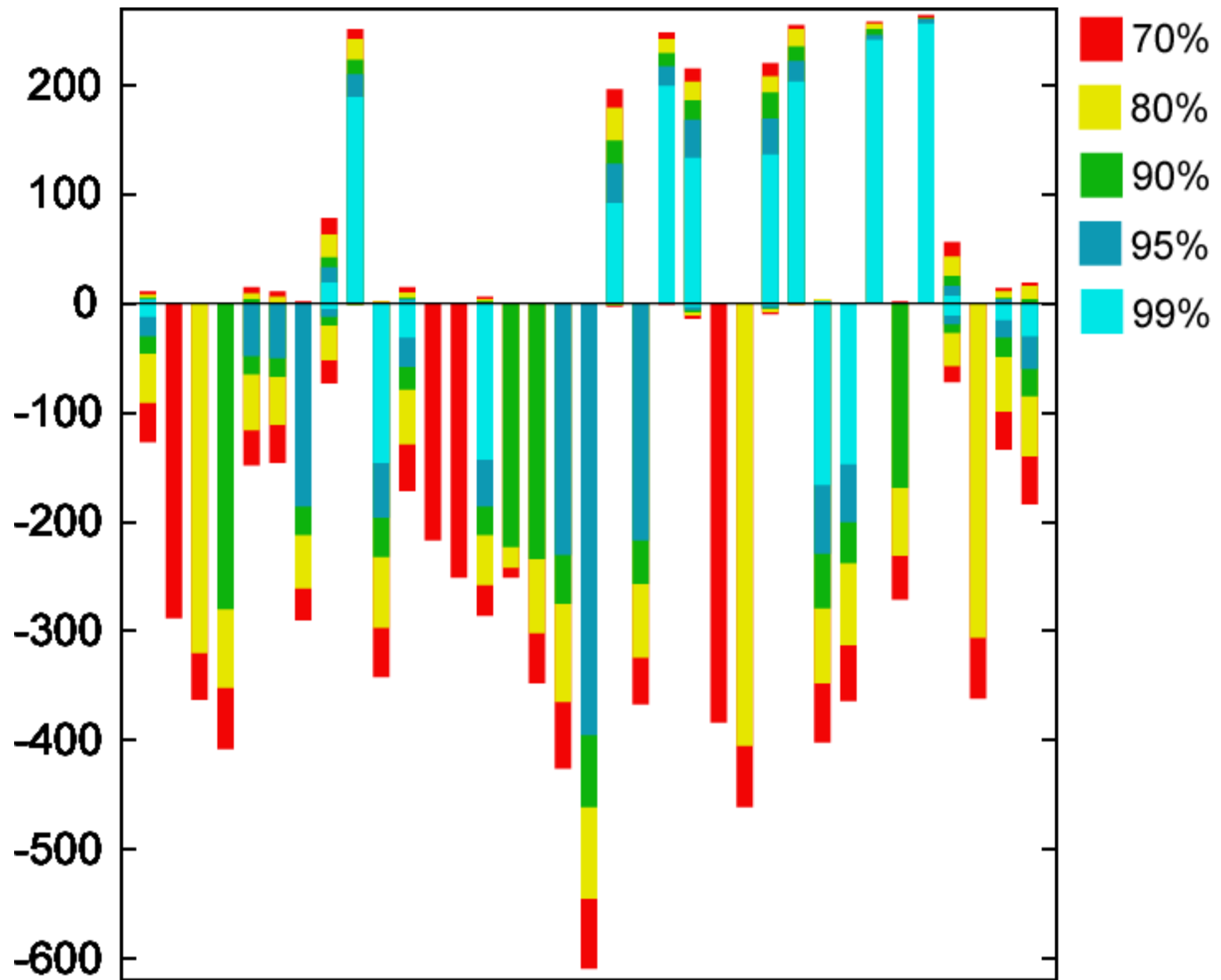


# “Lento”-plot of 35 supported bipartitions (out of 4082 possible)

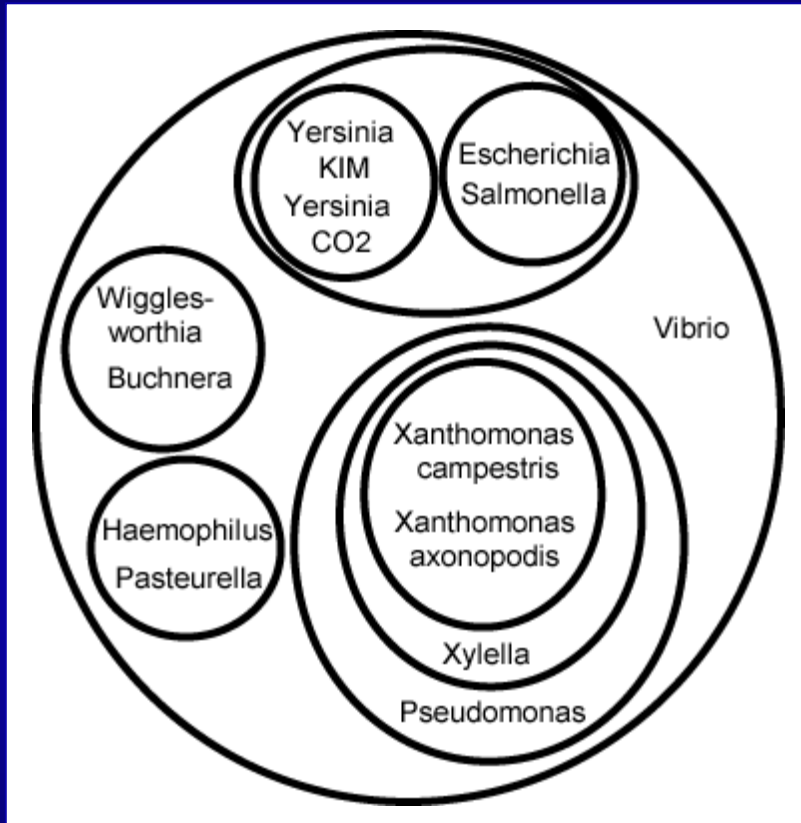
13  
gamma-  
proteobacterial  
genomes  
(258 putative  
orthologs):

- E.coli
- Buchnera
- Haemophilus
- Pasteurella
- Salmonella
- Yersinia pestis  
(2 strains)
- Vibrio
- Xanthomonas  
(2 sp.)
- Pseudomonas
- Wigglesworthia

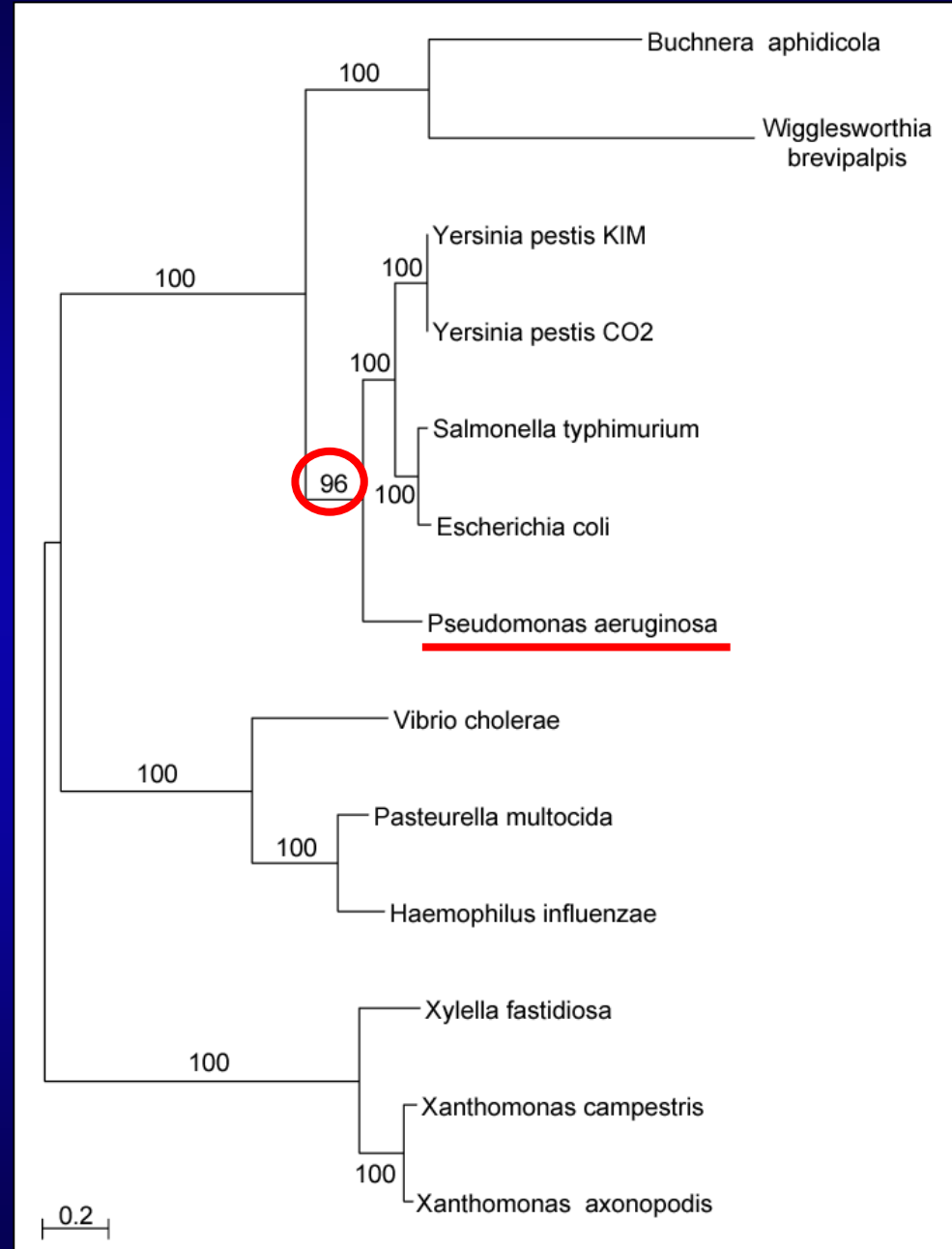
There are  
13,749,310,575  
possible  
unrooted tree  
topologies for  
13 genomes



# Consensus cluster of significantly supported bipartitions



# Phylogeny of virulence factor homologs (mviN)



# Case of Cyanobacteria

Based on 16S rRNA:

- 13 gamma proteobacteria have up to 19.8% sequence divergence,
- 10 cyanobacteria are at most 14% divergent.

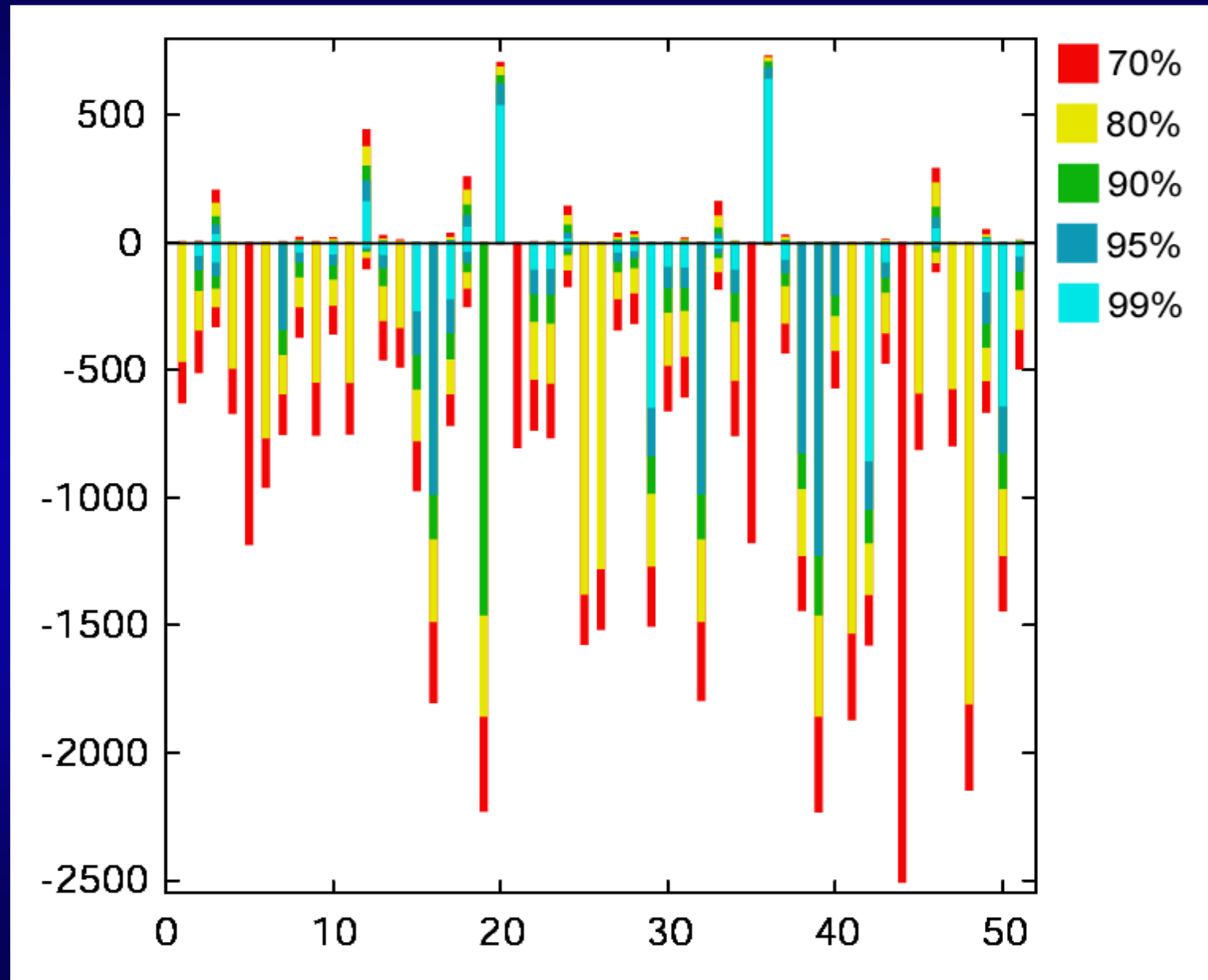


There are 678 orthologous genes detected by the reciprocal hit scheme.

# “Lento”-plot of 51 supported bipartitions (out of 501 possible)

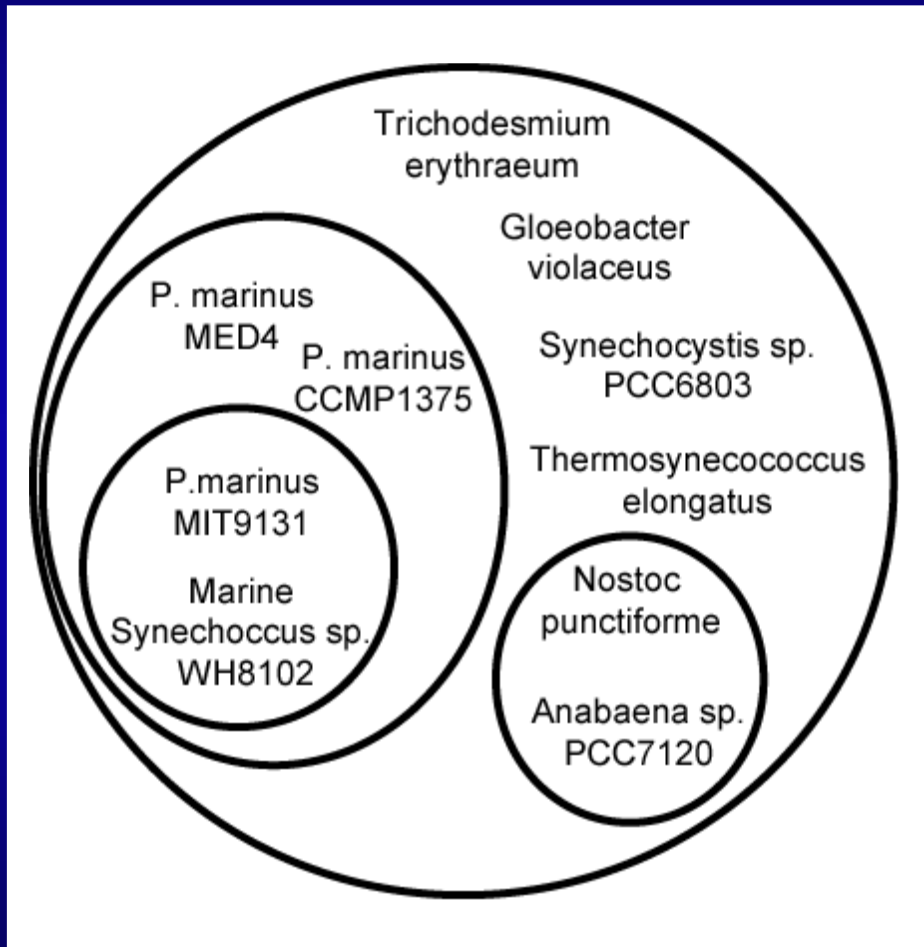
## 10 cyanobacteria:

- *Anabaena*
- *Trichodesmium*
- *Synechocystis* sp.
- *Prochlorococcus marinus* (3 strains)
- Marine *Synechococcus*
- *Thermosynechococcus elongatus*
- *Gloeobacter*
- *Nostoc punctioforme*

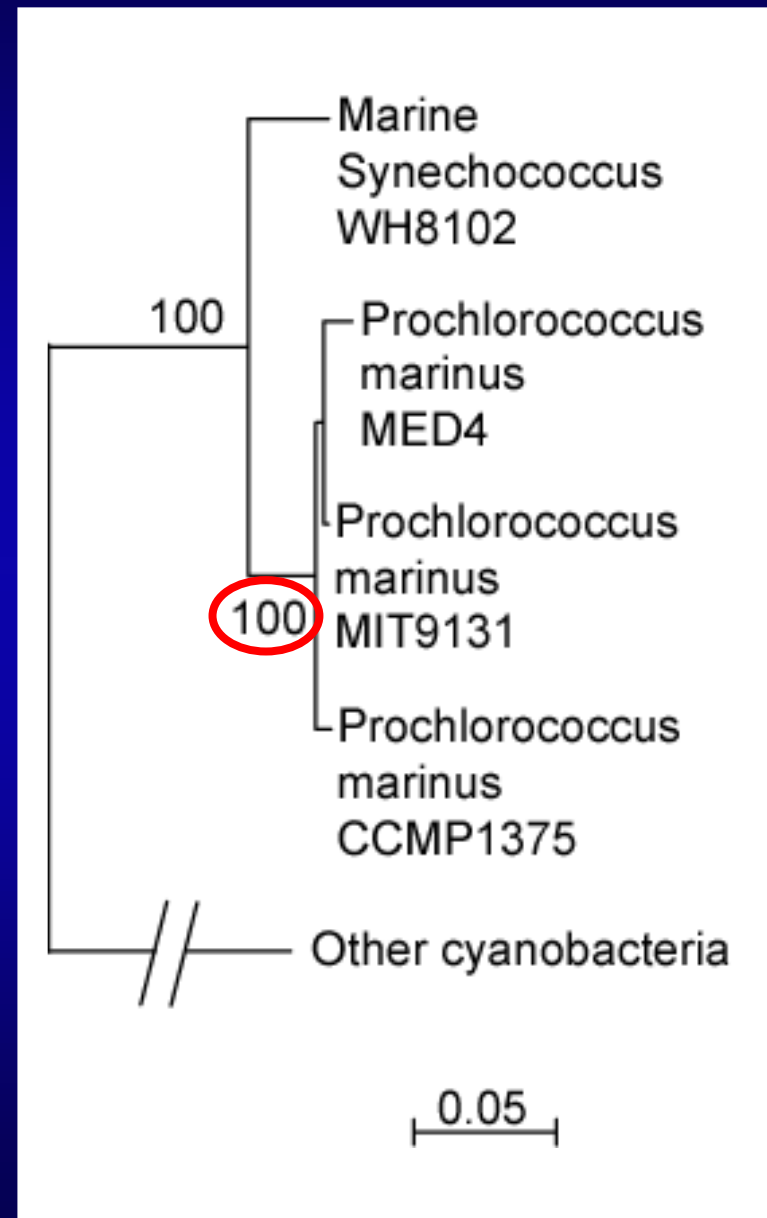




## Consensus cluster of significantly supported bipartitions



## The phylogeny of ribulose biphosphate carboxylase large subunit



## Other genes in conflict with the consensus at $\geq 99\%$ bootstrap support:

- ❖ cell division protein FtsH,
- ❖ translation initiation factor IF-2,
- ❖ ferredoxin, *petF*
- ❖ geranylgeranyl hydrogenase, *chlP*
- ❖ amidophosphoribosyltransferase,
- ❖ photosystem II reaction center core protein D2, *psbD*
- ❖ photosystem II CP43 core antenna protein, *psbC*
- ❖ photosystem II CP47 core antenna protein, *psbB*
- ❖ photosystem I reaction center core protein A2, *psaB*
- ❖ photosystem I reaction center core protein A1, *psaA*
- ❖ photosystem II manganese-stabilizing protein, *psbO*
- ❖ 5'-methylthioadenosine phosphorylase.

# Transfer of photosynthesis genes to and from *Prochlorococcus* viruses

Debbie Lindell<sup>†</sup>, Matthew B. Sullivan<sup>†‡</sup>, Zackary I. Johnson<sup>\*</sup>, Andrew C. Tolonen<sup>‡</sup>, Forest Rohwer<sup>§</sup>, and Sallie W. Chisholm<sup>\*¶</sup>

[www.pnas.org/cgi/dol/10.1073/pnas.0401526101](http://www.pnas.org/cgi/dol/10.1073/pnas.0401526101)

PNAS | July 27, 2004 | vol. 101 | no. 30 | 11013–11018

## Photosynthetic genes found in *Prochlorococcus* phages:

- PSII core reaction center protein D1 (*psbA*)
- PSII core reaction center protein D2 (*psbD*)
- ferredoxin (*petF*)
- plastocyanin (*petE*)
- HLIP cluster 14-type protein (*hli14* – high light inducible protein)

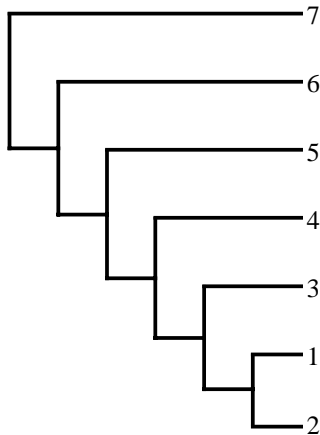
# CONCLUSIONS I

- Genomes are mosaic
- Support value mapping is a useful tool to dissect mosaic genomes
- While ML mapping can provide a quick assessment of genome mosaicism, it grossly overestimates reliability
- Analyzing extended datasets using embedded subtrees solves the problems associated with taxon sampling without sacrificing the visually appealing graphical representation

## CONCLUSIONS II

- Bipartition plots are a useful tool for comparative genome analyses. They allow to identify the plurality consensus cluster of genes contained in genomes as well as genes that conflict with the plurality consensus.
- In many instances majority or at least plurality signals are obtained from the analysis of individual genes.
- Sometimes clade-defining characteristics are among the genes that are transferred. E.g., for photosynthetic bacteria: plurality consensus phylogeny of genes  $\neq$  phylogeny of the chlorophyll biosynthetic enzymes.

$Q_1 = \{$   
 4 5 6 7  
 1 5 6 7  
 2 5 6 7  
 3 5 6 7  
 3 4 6 7  
 1 4 6 7  
 2 4 6 7  
 2 3 6 7  
 1 3 6 7  
 1 2 6 7  
 1 2 3 7  
 1 2 4 7  
 1 3 4 7  
 2 3 4 7  
 2 3 5 7  
 1 3 5 7  
 1 2 5 7  
 1 4 5 7  
 2 4 5 7  
 3 4 5 7  
 3 4 5 6  
 1 4 5 6  
 2 4 5 6  
 2 3 5 6  
 1 3 5 6  
 1 2 5 6  
 1 2 3 6  
 1 2 4 6  
 1 3 4 6  
 2 3 4 6  
 2 3 4 5  
 1 3 4 5  
 1 2 4 5  
 1 2 3 5  
 1 2 3 4  
 $\}$

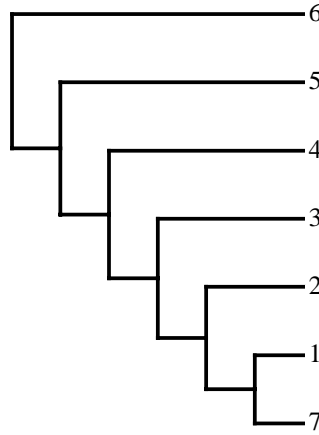


$B_1 = \{$   
 \*\* .....,  
 \*\*\* .....,  
 \*\*\*\* .....,  
 \*\*\*\*\* ..  
 $\}$

supported quartets  
 $Q_1 \cap Q_2 =$   
 $\{3 4 5 6, 1 4 5 6, 2 4 5 6, 2 3 5 6,$   
 $1 3 5 6, 1 2 5 6, 1 2 3 6, 1 2 4 6,$   
 $1 3 4 6, 2 3 4 6, 2 3 4 5, 1 3 4 5,$   
 $1 2 4 5, 1 2 3 5, 1 2 3 4\}$   
 supported bipartitions:  
 $B_1 \cap B_2 = \emptyset$

Illustration of a topology where  
 quartet analyses are more useful  
 than bipartition analyses

$Q_2 = \{$   
 3 4 5 6  
 1 4 5 6  
 7 4 5 6  
 2 4 5 6  
 2 3 5 6  
 1 3 5 6  
 7 3 5 6  
 7 2 5 6  
 1 2 5 6  
 1 7 5 6  
 1 7 2 6  
 1 7 3 6  
 1 2 3 6  
 7 2 3 6  
 7 2 4 6  
 1 2 4 6  
 1 7 4 6  
 1 3 4 6  
 7 3 4 6  
 2 3 4 6  
 2 3 4 5  
 1 3 4 5  
 7 3 4 5  
 7 2 4 5  
 1 2 4 5  
 1 7 4 5  
 1 7 2 5  
 1 7 3 5  
 1 2 3 5  
 7 2 3 5  
 7 2 3 4  
 1 2 3 4  
 1 7 3 4  
 1 7 2 4  
 1 7 2 3  
 $\}$



$B_2 = \{$   
 \* .....,  
 \*\* .....,  
 \*\*\* .....,  
 \*\*\*\* ..\*,  
 $\}$

# FUTURE RESEARCH

“Replace” bipartitions with  
 Embedded Quartets in  
 spectral analyses

+ Gene families that are not  
 represented in all genomes  
 can be included

+ adding more sequences  
 does not deteriorate support  
 values

+ a single “rogue” sequence  
 does not erase all of the  
 captured phylogenetic  
 information