



# Assessing the effects of compositional heterogeneity on phylogenomic analyses

---

Denis Baurain<sup>1,2,3</sup>, Robert G. Beiko<sup>2</sup>, and Mark A. Ragan<sup>2</sup>

<sup>1</sup> Université de Liège / FNRS, Belgium

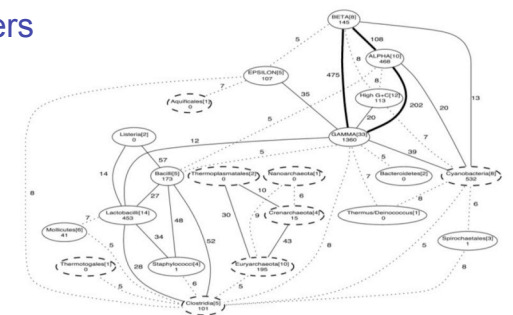
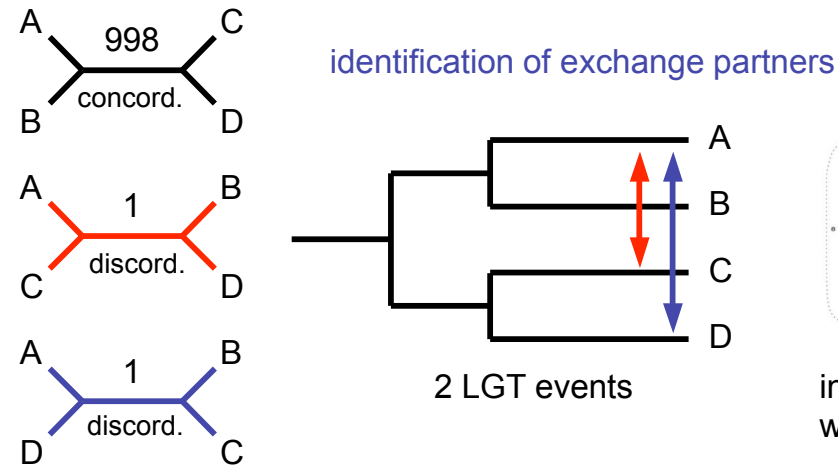
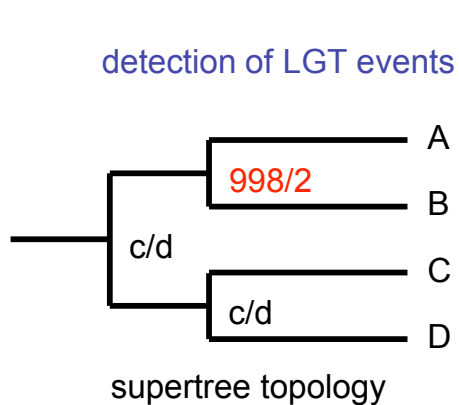
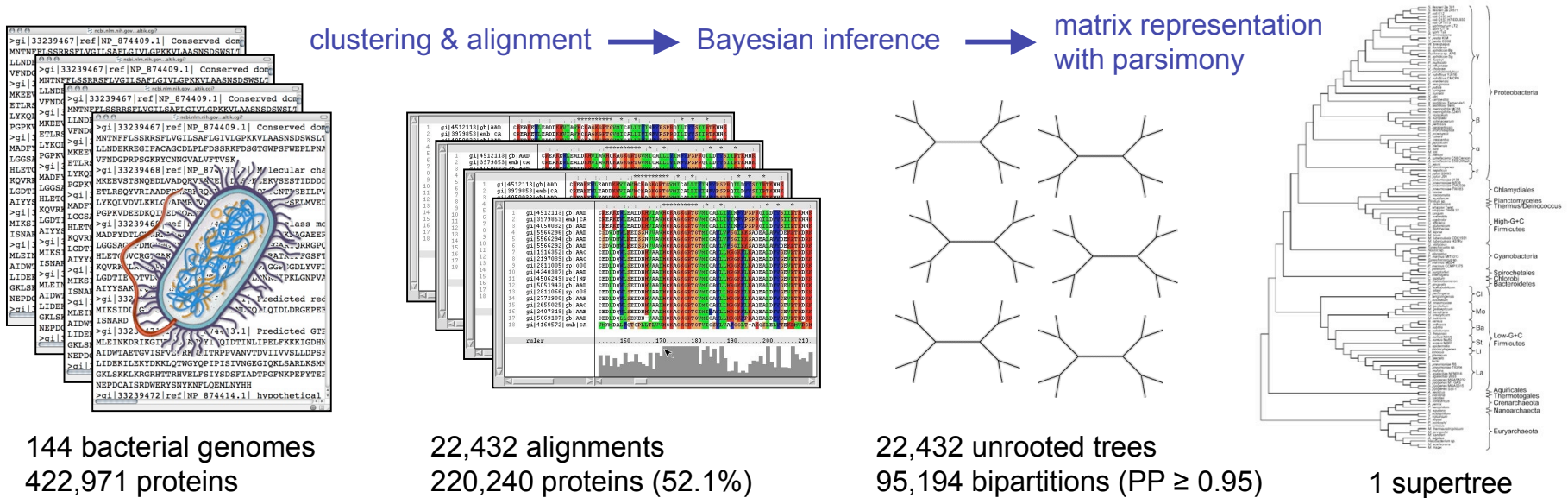
<sup>2</sup> University of Queensland / ARC Centre in Bioinformatics, Australia

<sup>3</sup> Université de Montréal

November 4<sup>th</sup>, 2005

# Background

Beiko *et al.* (2005) PNAS **102**:14332-14337

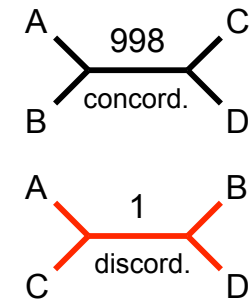


inheritance pattern largely vertical with 'highways' of gene sharing

# LGT or methodological issues?

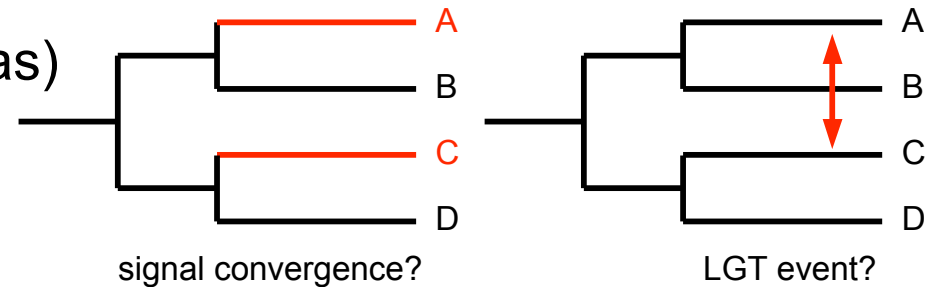
- tests for **systematic biases**

- clustering strategy (e.g. cluster size)
- alignment quality (e.g. sequence length variation)
- phylogenetic inference (e.g. alignment size)



- **conflicting signals in the data**

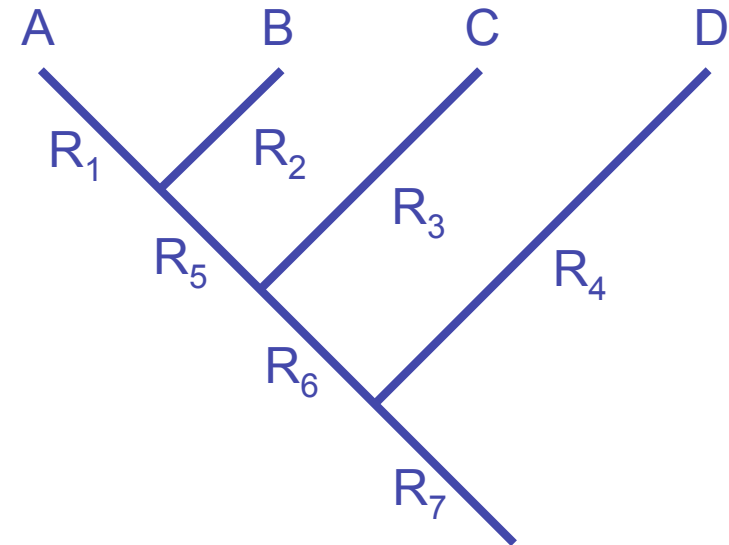
- history
- composition (e.g. GC bias)
- rate (e.g. LBA)
- other signals



- « Any signal having experienced **convergence in nonsister lineages** will affect recovery of the historical signal. »

# Phylogenetic assumptions

- sites evolve **independently** and identically using a Markov process given by  $R$  (e.g. GTR model)
- for practical reasons, we assume that the sites have evolved under **reversible, homogeneous** and **stationary** conditions
- **note** – new models allow to safely violate some of these assumptions



$$R_k = \begin{bmatrix} -\sum_{j \neq A} \pi_j \alpha_{Aj} & \pi_C \alpha_{AC} & \pi_G \alpha_{AG} & \pi_T \alpha_{AT} \\ \pi_A \alpha_{AC} & -\sum_{j \neq C} \pi_j \alpha_{Cj} & \pi_G \alpha_{CG} & \pi_T \alpha_{CT} \\ \pi_A \alpha_{AG} & \pi_C \alpha_{CG} & -\sum_{j \neq G} \pi_j \alpha_{Gj} & \pi_T \alpha_{GT} \\ \pi_A \alpha_{AT} & \pi_C \alpha_{CT} & \pi_G \alpha_{GT} & -\sum_{j \neq T} \pi_j \alpha_{Tj} \end{bmatrix}$$

$\pi_A, \dots, \pi_T$  – nucleotide frequencies

$\alpha_{kij}$  – conditional rates of change

$\alpha_{kij} = \alpha_{kji}$  – reversibility

$R_1 = R_2 = \dots = R_6 = R_7$  – homogeneity

$\pi_{kj} = f_{0j}$  – stationarity

# Statistical tests for stationarity

S<sub>1</sub> ATGGTACAATGCGGCATGTACTCGCGATATCGACGATACG

S<sub>2</sub> ATCGAACGATGTGGCGTACACTCACGTTACCGACACGACG

- matched-pair tests of homogeneity
  - Bowker's (1948) test for *symmetry*

$$\chi_{bowker}^2 = \sum_{i < j} \frac{(x_{ij} - x_{ji})^2}{x_{ij} + x_{ji}}$$

- Stuart's (1955) test for *marginal homogeneity*

$$\chi_{stuart}^2 = d' S^{-1} d \quad \text{with } d = x_{i.} - x_{.i}$$

$$\text{and } s_{ii} = x_{i.} + x_{.i} - 2x_{ii}, s_{ij} = -(x_{ij} + x_{ji}), i \neq j$$

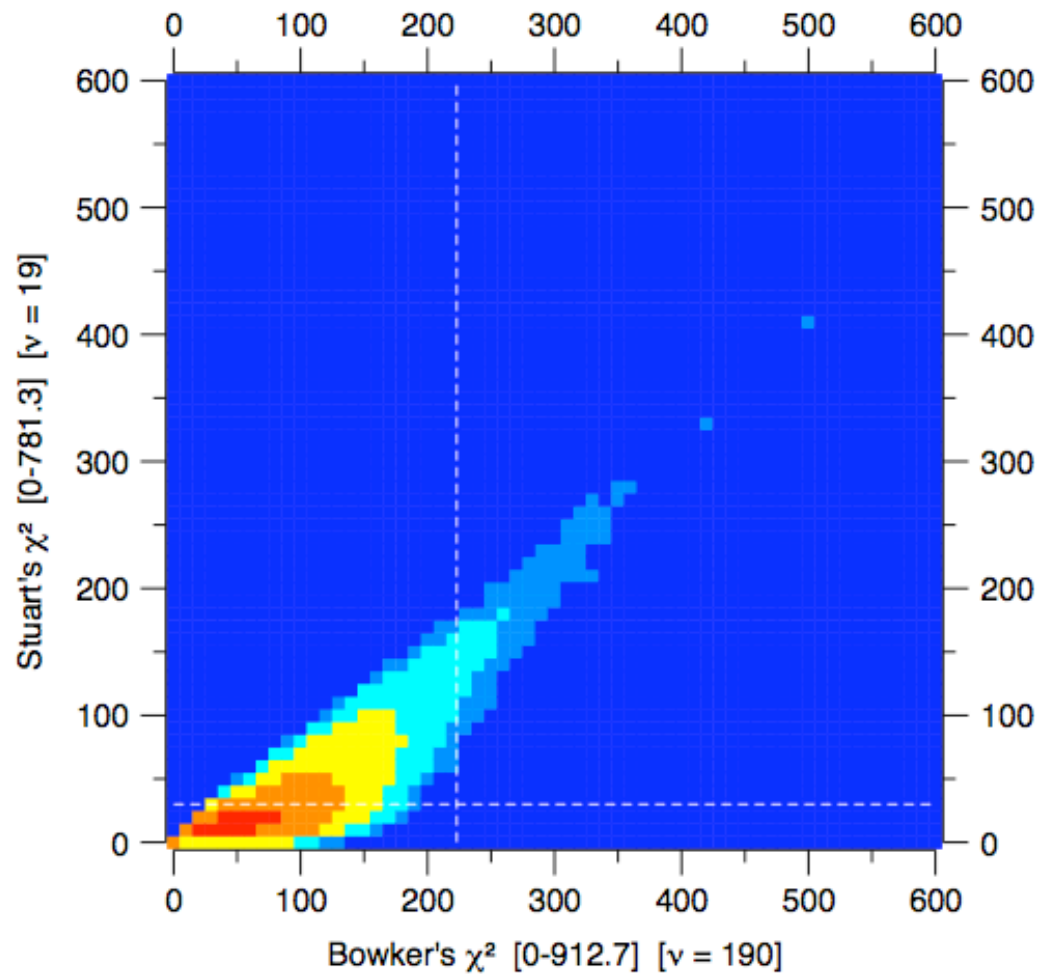
- 'traditional' homogeneity tests
  - compositional  $\chi^2$

$$\chi_{compos.}^2 = \sum_{i=1}^K \frac{(x_{i.} - e_i)^2}{e_i} \quad \text{with } e_i = \frac{x_{i.} + x_{.i}}{2}$$

	A	C	G	T	Σ
A	7	0	3	1	11
C	1	8	1	2	12
G	2	0	7	1	10
T	1	1	0	5	7
Σ	11	9	11	9	40

two-way contingency table

# Correlations among all three tests



## failure statistics

[AA, n = 2815041, p = 0.05]

Bowker's $\chi^2$	0.23%
Stuart's $\chi^2$	27.53%
compositional $\chi^2$	2.47%

Pearson's correlation	Bowker's $\chi^2$	Stuart's $\chi^2$	compositional $\chi^2$
Bowker's $\chi^2$	-	0.752	0.712
Stuart's $\chi^2$	0.752	-	0.946
compositional $\chi^2$	0.712	0.946	-



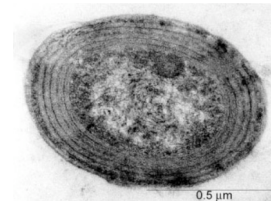
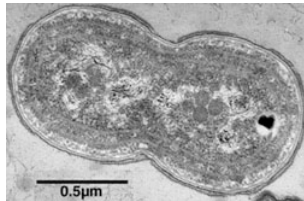
# Ranking of the worst players

rank	fail. (#)	pairs (#)	fail. (%)	organism a	organism b
1	317	1274	24.88	<i>Synechococcus</i> WH8102	<i>Prochlorococcus</i> MED4
2	269	1245	21.61	<i>Prochlorococcus</i> MIT9313	<i>Prochlorococcus</i> MED4
3	204	835	24.43	<i>Thermosynechococcus</i> BP-1	<i>Prochlorococcus</i> MED4
4	203	749	27.10	<i>Gloeobacter violaceus</i>	<i>Prochlorococcus</i> MED4
5	168	868	19.35	<i>Synechocystis</i> PCC6803	<i>Prochlorococcus</i> MED4
6	160	390	41.03	<i>Wigglesworthia brevipalpis</i>	<i>Pseudomonas</i> PA01
...	...	...	...	...	...
12	152	1301	11.68	<i>Synechococcus</i> WH8102	<i>Prochlorococcus</i> SS120
...	...	...	...	...	...
29	131	774	16.93	<i>Gloeobacter violaceus</i>	<i>Prochlorococcus</i> SS120
...	...	...	...	...	...

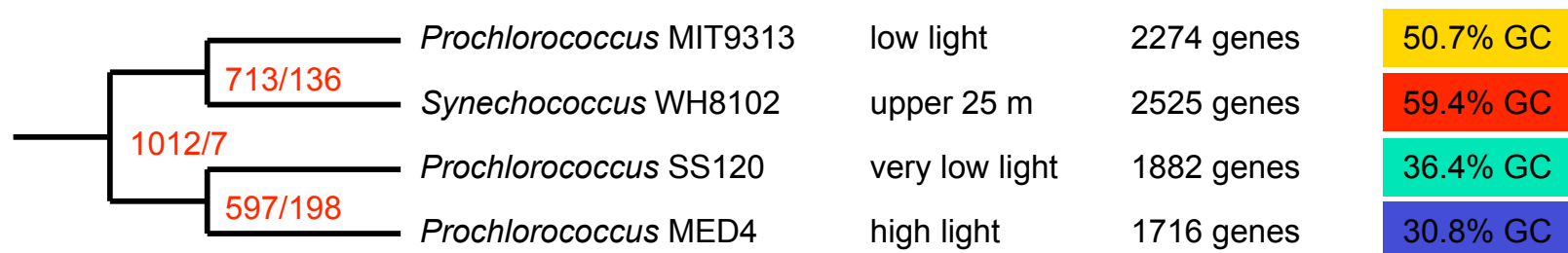
[Bowker's  $\chi^2$ , AA, n = 2815041, 10296 organism pairs]

# Who are the picocyanobacteria?

- both *Prochlorococcus* and *Synechococcus* are part of the **picophytoplankton** (tiny organisms: cell  $\emptyset < 1 \mu\text{m}$ )
  - they account for 1/3 of Earth's primary biomass production
  - all known members of this group are **96% similar at the rRNA level** (but have quite different gene contents)
  - *Synechococcus* found 25 ya / *Prochlorococcus* found 15 ya

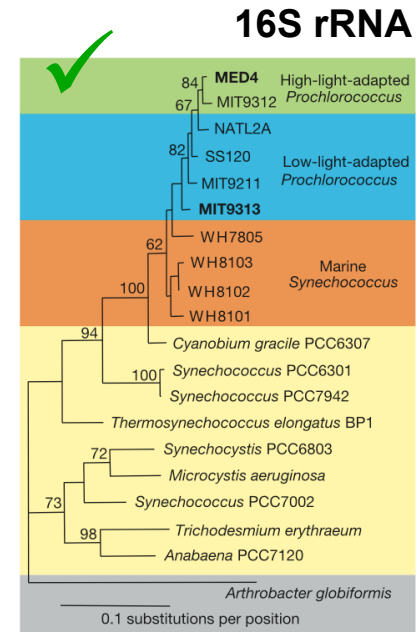
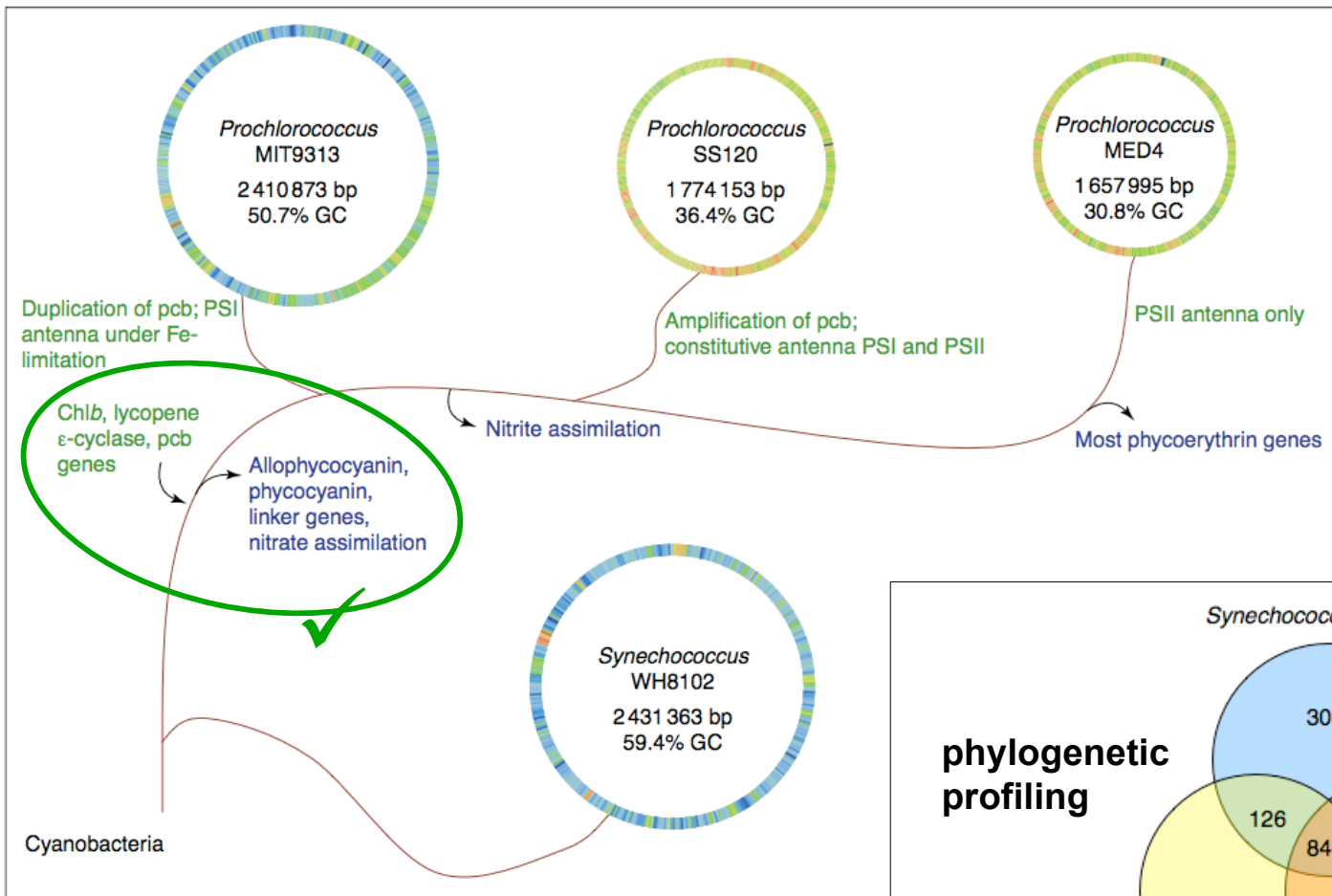


- we had 4 genomes in our 144-species dataset





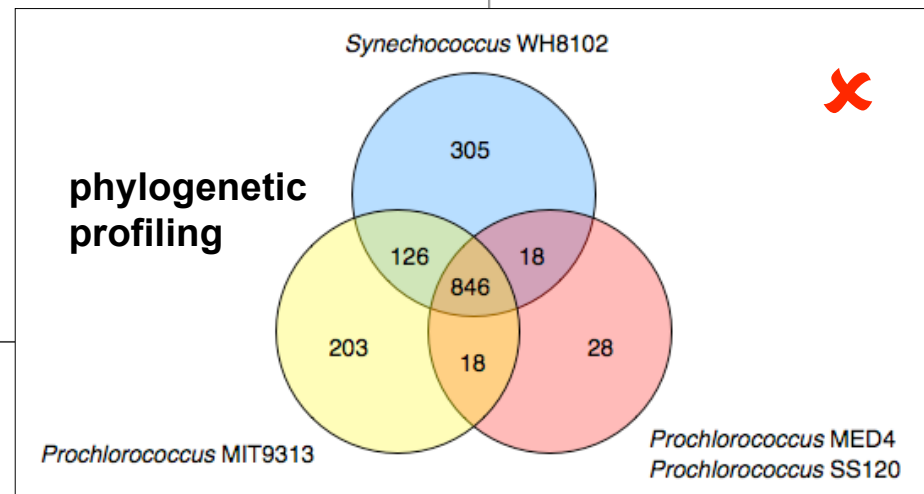
# Is *Prochlorococcus* monophyletic?



Rocap et al. (2003) Nature 424:1042-1047

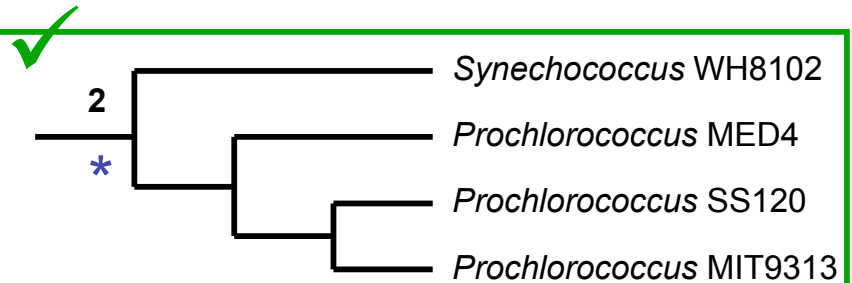
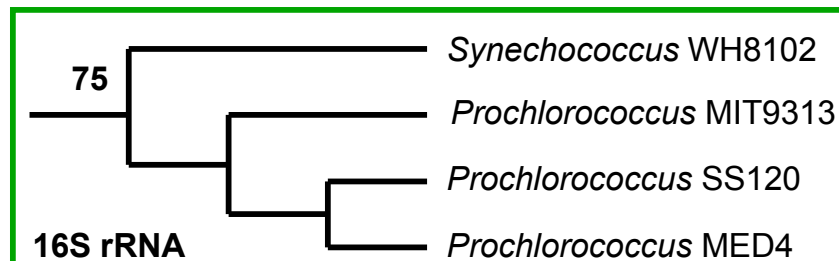
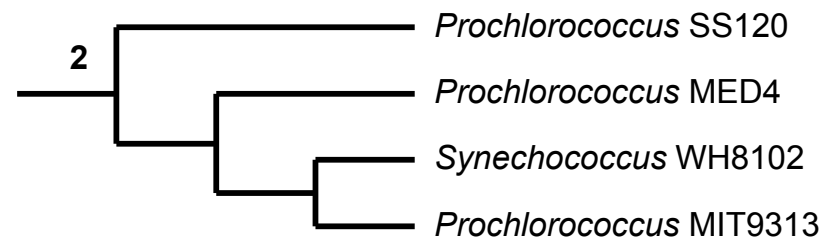
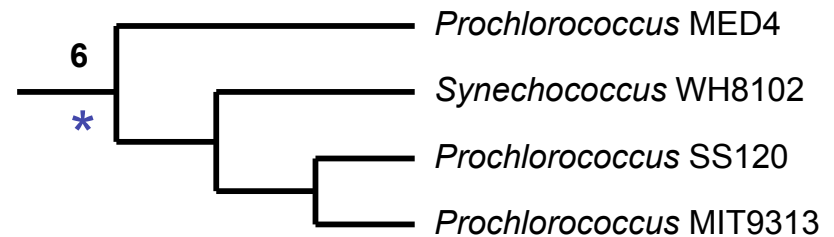
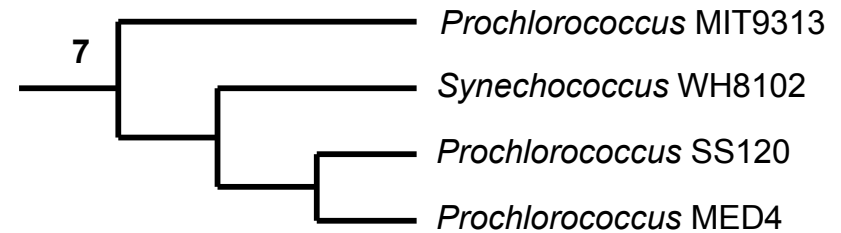
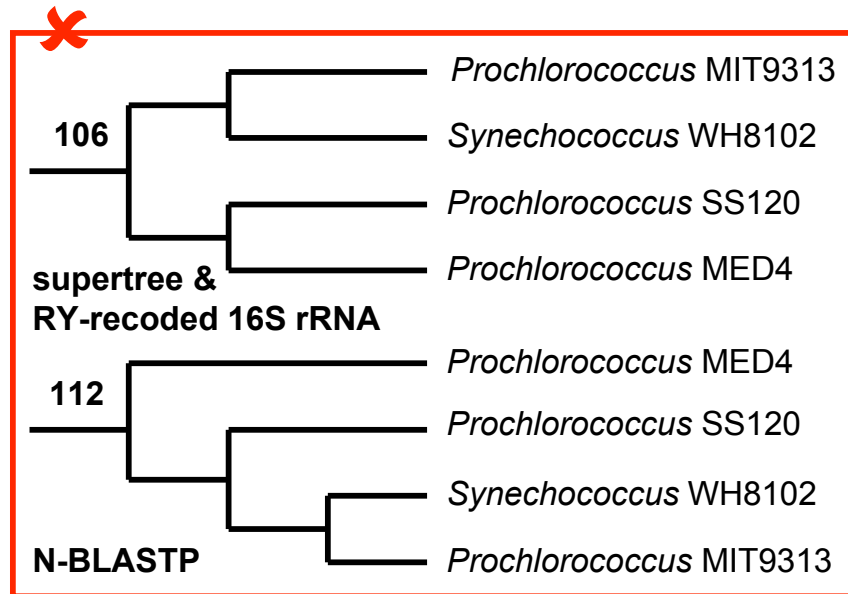
## photosynthetic apparatus

Hess (2004) Curr Opin Biotech 15:191-198



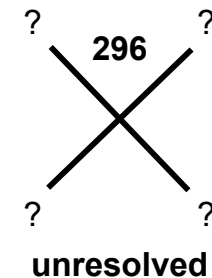
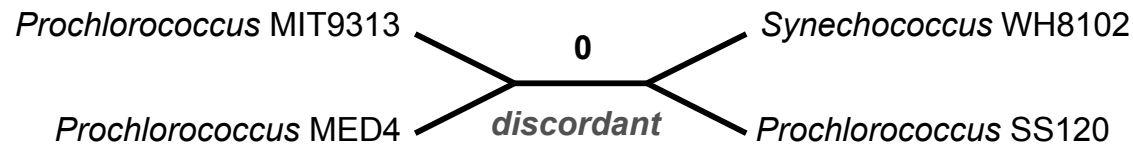
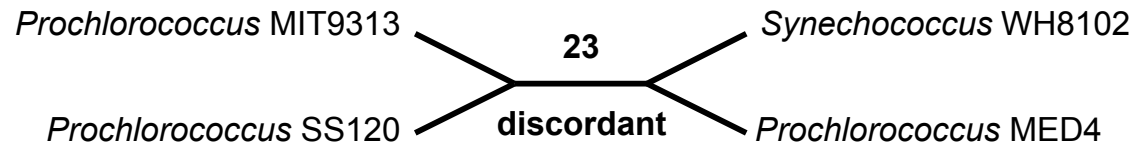
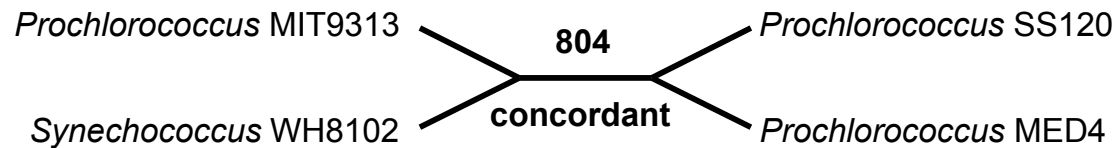
# 7 rooted trees...

alignments of [size  $\geq 6$ ] for which the monophyly of the four picocyanobacteria is highly supported (n = 819)  
 only fully resolved topologies are considered (n = 310)  
 blue stars (\*) denote *really* discordant topologies



# ... fold to 2 unrooted topologies

same dataset + all alignments of  $[4 \leq \text{size} \leq 5]$  that include the four picocyanobacteria ( $n = 304$ )  
all topologies are considered ( $n = 1123$ )



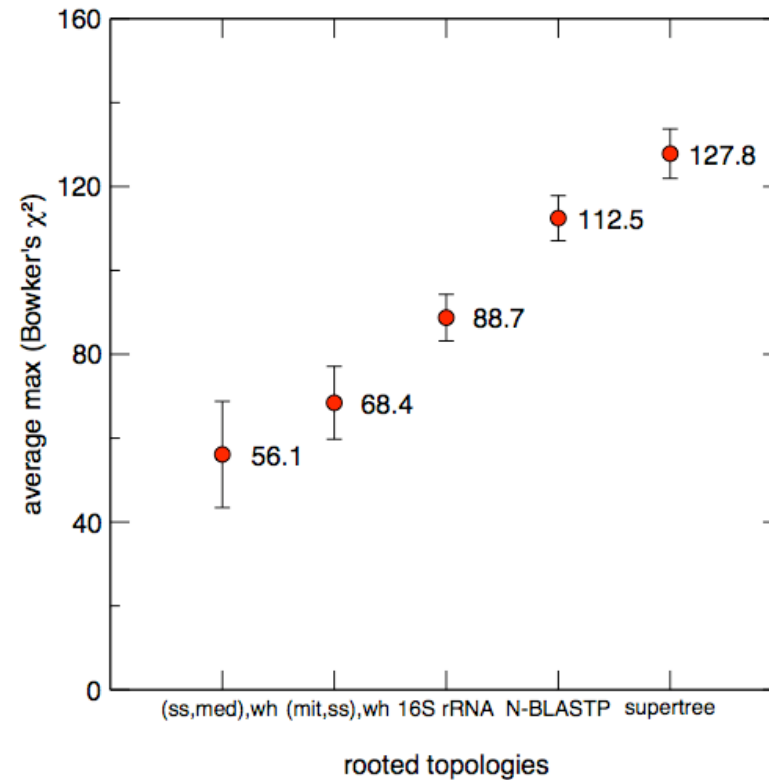
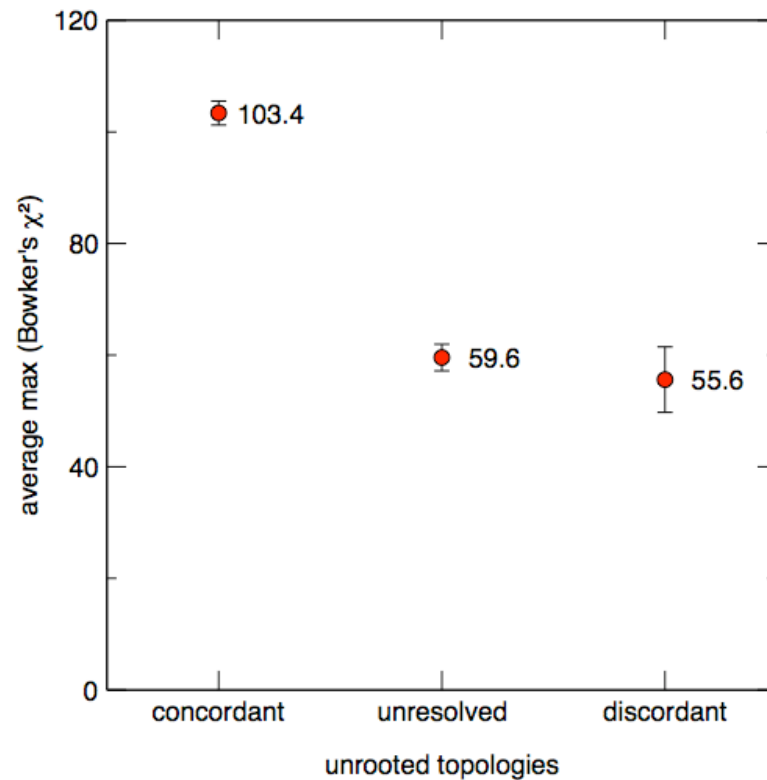
having binned all alignments leading to one given topology, we look for features of the compositional signal (compositional bias) that would be characteristic of that topology

# Signals that work together

$$\chi_{bowker}^2 = \sum_{i < j} \frac{(x_{ij} - x_{ji})^2}{x_{ij} + x_{ji}}$$

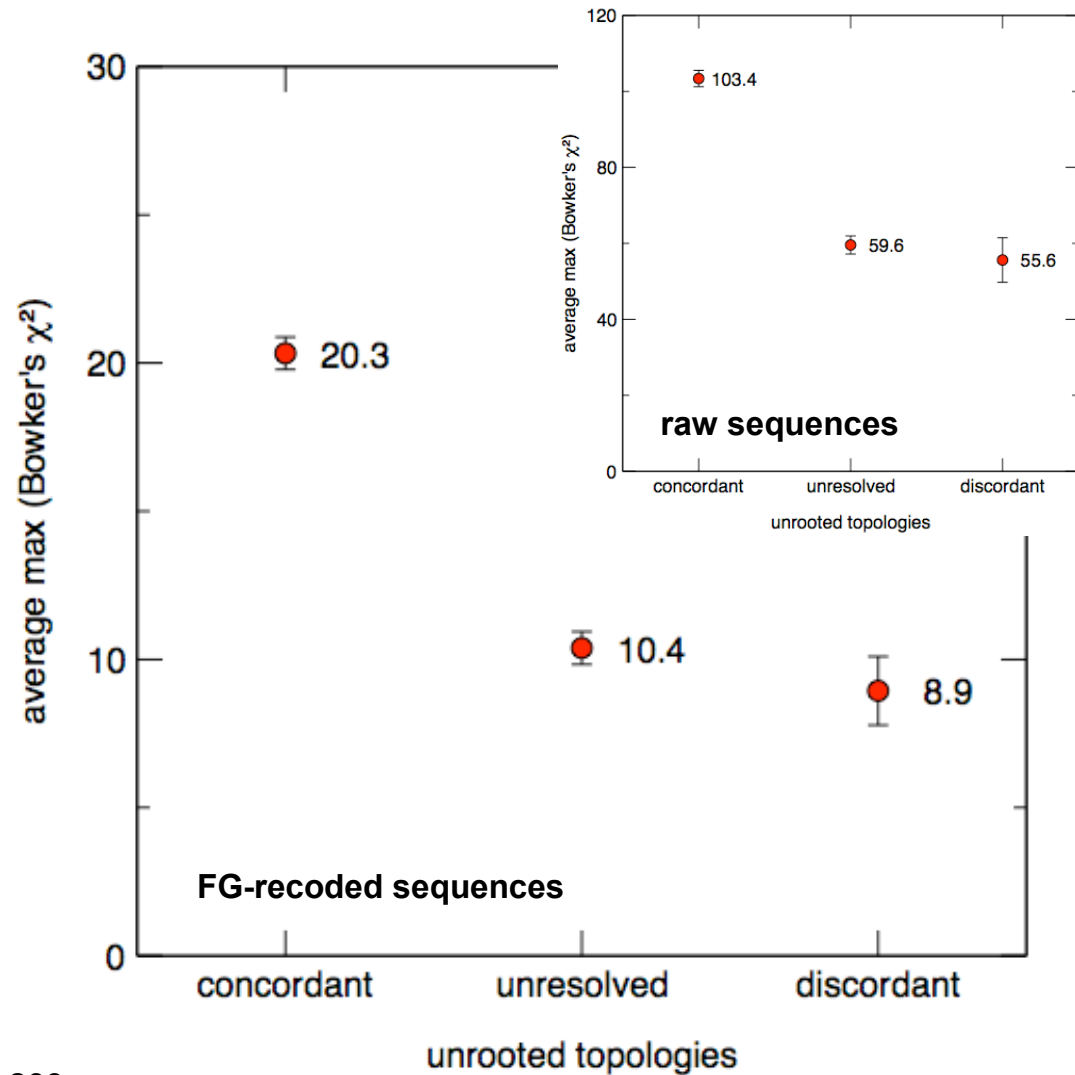
only bins of [size  $\geq 6$ ] are considered  
bars denote standard error

topologies are sorted on y values

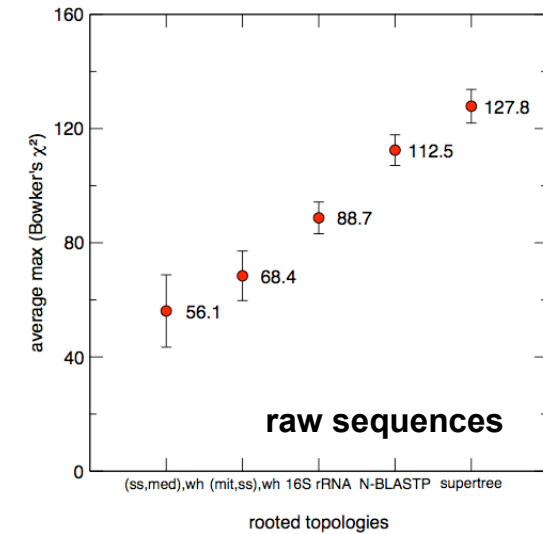
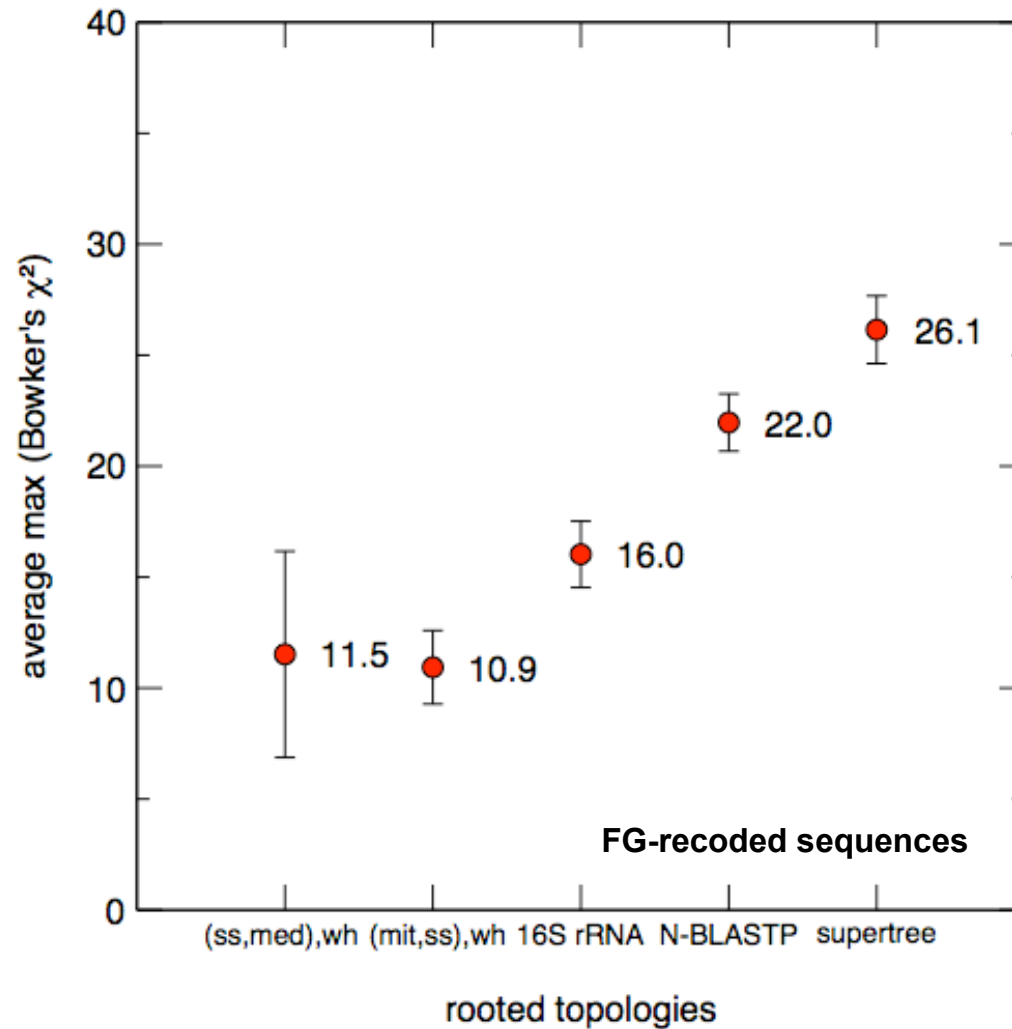


# GC bias at the protein level?

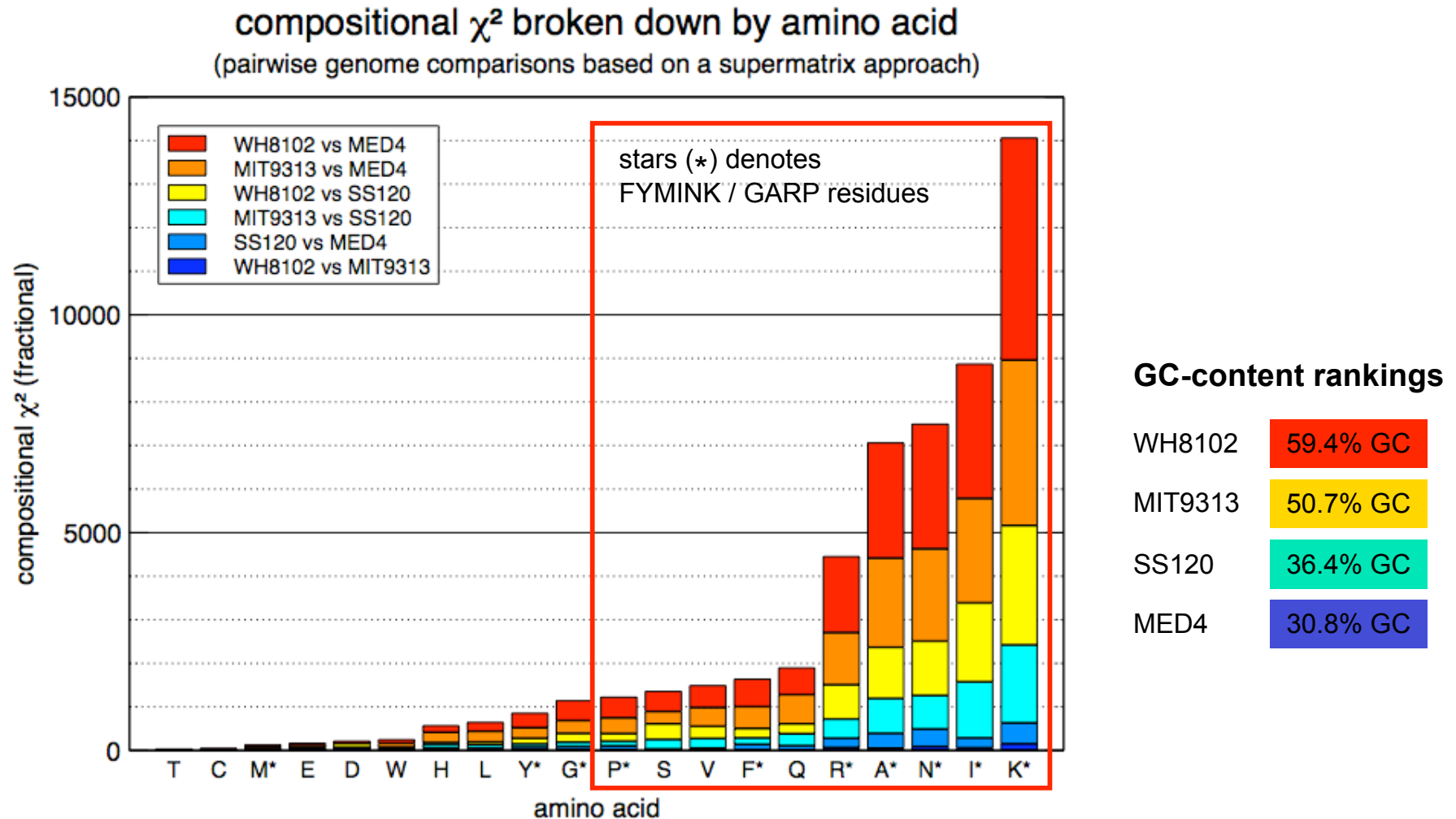
		2nd position AT		2nd position GC	
1st position AT	aar	atr	agr	acn	
	aay	aty	agy	acn	
	tar	ttr	tgr	tcn	
	tay	tty	tgy	tcn	
1st position GC	gar	gtn	ggn	gcn	
	gay	gtn	ggn	gcn	
	car	ctn	cgn	ccn	
	cay	ctn	cgn	ccn	
		K M		S T	
		N I		W S	
		L		C S	
		Y F			
		E V		G A	
		D		R P	
		Q L			
		H			



# GC bias at the protein level? (2)



# The case for a closer look



'supergene' made of 1,485 genes found in at least two (out of four) organisms (469,682 positions; 15% missing)



# Conclusions

---

- the compositional heterogeneity definitely **has an impact** on phylogenomic analyses
  - here, the compositional signal likely exaggerates the historical signal (i.e. non-monophyly of *Prochlorococcus*)
  - in other circumstances, it could be the opposite
- while the compositional heterogeneity is obvious at both the DNA and protein levels, the propagation of the **GC bias** is not limited to **FYMINK / GARP** codons
- in their ‘holy war’ against systematic biases, phylogenomic analyses would certainly benefit from **newer evolutionary models** that can deal with the violation of the stationarity assumption
  - e.g. Galtier and Gouy (1998) Mol Biol Evol **15**:871-879



# Acknowledgments

- data, help and advice
  - Robert Charlebois (NeuroGadgets Inc.)
  - Nick Hamilton (UQ)
  - Tim Harlow (UQ / ACB)
  - Lars S. Jermiin (USyd)
- funding
  - FNRS / ULg (Belgium)
  - ACB (Australia)



**IMB** Institute for Molecular Bioscience



Fonds National de la Recherche Scientifique



Australian Research Council Centre in Bioinformatics