

Diagnosing phylogenetic conflict between genes and proteins: Evidence for the origin of plastids

Blaise Li Peter G. Foster T. Martin Embley Cymon J. Cox

SMBE Meeting - Dublin 2012 - 24/06/2012

The questions

- ▶ Archaeplastida (plants) = Rhodophyta (red algae) + Glaucophyta + Viridiplantae (green plants)

The questions

- ▶ Archaeplastida (plants) = Rhodophyta (red algae) + Glaucophyta + Viridiplantae (green plants)
→ Characterised by a plastid, with endosymbiotic origin within cyanobacteria

¹Notation adapted from Criscuolo and Gribaldo (2011)

The questions

- ▶ Archaeplastida (plants) = **Rhodophyta** (red algae) + **Glaucophyta** + **Viridiplantae** (green plants)
→ Characterised by a plastid, with endosymbiotic origin within **cyanobacteria**
- ▶ Groups of cyanobacteria already identified: **NOST-1, OSC-2, SPM-3, SO-6, GBACT, UNIT+¹**

¹Notation adapted from Criscuolo and Gribaldo (2011)

The questions

- ▶ Archaeplastida (plants) = **Rhodophyta** (red algae) + **Glaucophyta** + **Viridiplantae** (green plants)
→ Characterised by a plastid, with endosymbiotic origin within **cyanobacteria**
- ▶ Groups of cyanobacteria already identified: **NOST-1, OSC-2, SPM-3, SO-6, GBACT, UNIT+¹**
- ▶ What are the relationships between plastids and cyanobacteria?

¹Notation adapted from Criscuolo and Gribaldo (2011)

The questions

- ▶ Archaeplastida (plants) = **Rhodophyta** (red algae) + **Glaucophyta** + **Viridiplantae** (green plants)
→ Characterised by a plastid, with endosymbiotic origin within **cyanobacteria**
- ▶ Groups of cyanobacteria already identified: **NOST-1, OSC-2, SPM-3, SO-6, GBACT, UNIT+¹**
- ▶ What are the relationships between plastids and cyanobacteria?
- ▶ Why do the gene tree and the protein tree give a different answer?

¹Notation adapted from Criscuolo and Gribaldo (2011)

Dataset

- ▶ 42 taxa, including 8 outgroup (non-cyano)bacteria, 16 cyanobacteria, and plastids from 1 Glaucophyta, 4 Rhodophyta and 13 Viridiplantae

Dataset

- ▶ 42 taxa, including 8 outgroup (non-cyano)bacteria, 16 cyanobacteria, and plastids from 1 Glaucophyta, 4 Rhodophyta and 13 Viridiplantae
- ▶ 75 protein-coding genes, but 452 missing sequences (i.e. 14% overall, and up to 38 genes missing for one of the outgroup taxa)

Dataset

- ▶ 42 taxa, including 8 outgroup (non-cyano)bacteria, 16 cyanobacteria, and plastids from 1 Glaucophyta, 4 Rhodophyta and 13 Viridiplantae
- ▶ 75 protein-coding genes, but 452 missing sequences (i.e. 14% overall, and up to 38 genes missing for one of the outgroup taxa)
- ▶ Concatenated dataset (cg75)

Analysis method

- ▶ Nucleotide dataset (cg75)

Analysis method

- ▶ Nucleotide dataset (cg75) and amino-acid translation (cp75)

Analysis method

- ▶ Nucleotide dataset (cg75) and amino-acid translation (cp75)
- ▶ Various recodings of cg75 using IUPAC ambiguity codes to remove signal associated with some synonymous substitutions.

Analysis method

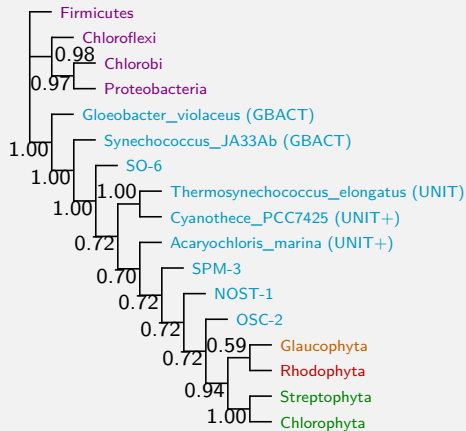
- ▶ Nucleotide dataset (cg75) and amino-acid translation (cp75)
- ▶ Various recodings of cg75 using IUPAC ambiguity codes to remove signal associated with some synonymous substitutions. e.g for His: CAC ↔ CAT (synonymous codons) → CAY (degenerate codon)

Analysis method

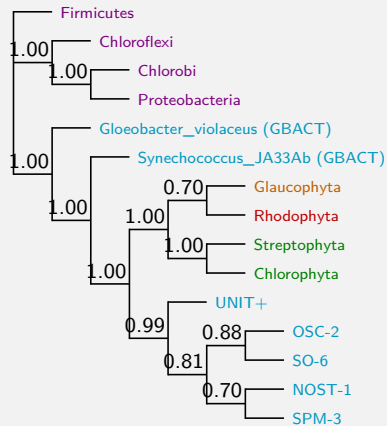
- ▶ Nucleotide dataset (cg75) and amino-acid translation (cp75)
- ▶ Various recodings of cg75 using IUPAC ambiguity codes to remove signal associated with some synonymous substitutions. e.g for His: CAC ↔ CAT (synonymous codons) → CAY (degenerate codon)
- ▶ Maximum likelihood bootstrap analyses (200 resamplings) using RAxML, with GTR+I+Γ (or LG+I+Γ for cp75, chosen using ProtTest)

Unrecoded ML bootstrap analyses

translation →



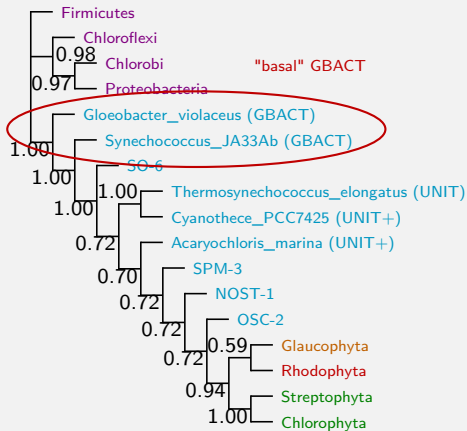
cg75



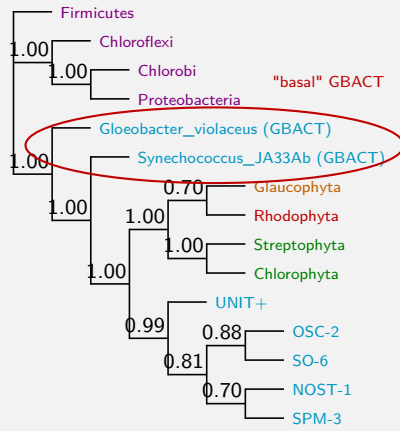
cp75

Unrecoded ML bootstrap analyses

translation →



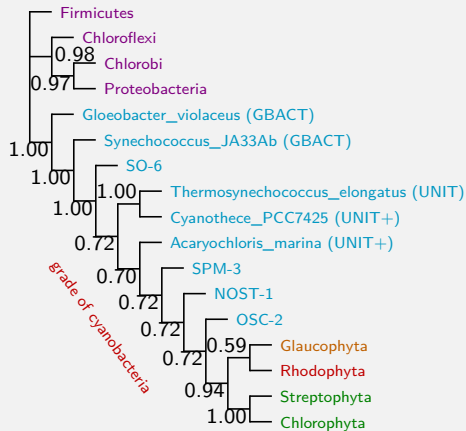
cg75



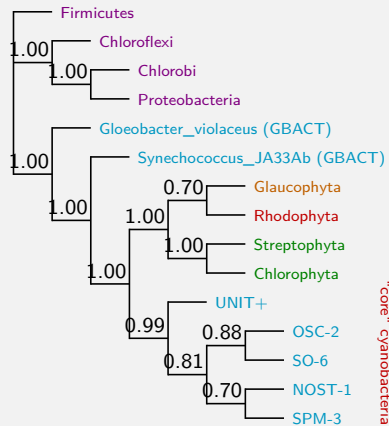
cp75

Unrecoded ML bootstrap analyses

translation →



cg75



cp75

Unrecoded ML bootstrap analyses

- ▶ cp75 is a direct translation of cg75

Unrecoded ML bootstrap analyses

- ▶ cp75 is a direct translation of cg75
→ The trees should be the same.

Unrecoded ML bootstrap analyses

- ▶ cp75 is a direct translation of cg75
→ The trees should be the same.
- ▶ But the analyses conflict in the identification of the plastid sister-group.

Unrecoded ML bootstrap analyses

- ▶ cp75 is a direct translation of cg75
→ The trees should be the same.
- ▶ But the analyses conflict in the identification of the plastid sister-group.
→ Something is not well modelled.

Unrecoded ML bootstrap analyses

- ▶ cp75 is a direct translation of cg75
→ The trees should be the same.
- ▶ But the analyses conflict in the identification of the plastid sister-group.
→ Something is not well modelled.
- ▶ Data recoding experiments to identify some causes of model mis-specification.

Degenerating synonymous 3rd codon positions

	T		C		A		G	
T	TTT	Phe	TCT	Ser	TAT	Tyr	TGT	Cys
	TTC		TCC		TAC		TGC	
	TTA	Leu	TCA		TAA	Ter	TGA	Ter
	TTG		TCG		TAG		TGG	Trp
C	CTT	Leu	CCT	Pro	CAT	His	CGT	Arg
	CTC		CCC		CAC		CGC	
	CTA		CCA		CAA	Gln	CGA	
	CTG		CCG		CAG		CGG	
A	ATT	Ile	ACT	Thr	AAT	Asn	AGT	Ser
	ATC		ACC		AAC		AGC	
	ATA		ACA		AAA	Lys	AGA	
	ATG	Met	ACG		AAG		AGG	Arg
G	GTT	Val	GCT	Ala	GAT	Asp	GGT	Gly
	GTC		GCC		GAC		GGC	
	GTA		GCA		GAA	Glu	GGA	
	GTG		GCG		GAG		GGG	

Example:

```

ATGAAACAGTTGCTGGAAGCCGGTGTTCACTTC
TTATCAGAACTTTTAGATGCCAGCGCTCACATA
TTAGAAGATATGATTCAAAGTGGAATGCATTTT
CTCGAACAAATGCTAGATGTAGGTGTTCAATTTT
TTACAGTCAATGCTTGAAGCTGGTGTTCACTTT
TTAGAAGCACTTTTAGAGGCTGGTGTTCATTTT
TTGGAAGAAATGATGGAAGCGGGTATTCATTTT
TTGGAACAAATGATGGAAGCAGGAGTCCATTTT
ATTCAGCAATTATTGGAAGCAGGAGTCCATCTG
ATTGAAGAATTGTTAGAAGCTGGCGTGCATTTT
TTAGAACAAATGTTAGATGCAGGTGTACATTTT
TTACAAAAAATGATTGAAGCTGGTGTTCATTTT
CTATCAGAAATGATGGAAGCTGGTGTTCATTTT
CTGCCGCAAATGCTGGAAGCCGGTGTCCATTTT
TTAGAAGAAATGATGGAAGCAGGGGTCCATTTT
TTAGCAGAATTACTAGAAGCGGGCGTTCAATTTT
  
```

Degenerating synonymous 3rd codon positions

	T		C		A		G	
T	TTT	Phe	TCT	Ser	TAT	Tyr	TGT	Cys
	TTC		TCC		TAC		TGC	
	TTA	Leu	TCA		TAA	Ter	TGA	Ter
	TTG		TCG		TAG		TGG	Trp
C	CTT	Leu	CCT	Pro	CAT	His	CGT	Arg
	CTC		CCC		CAC		CGC	
	CTA		CCA		CAA	Gln	CGA	
	CTG		CCG		CAG		CGG	
A	ATT	Ile	ACT	Thr	AAT	Asn	AGT	Ser
	ATC		ACC		AAC		AGC	
	ATA		ACA		AAA	Lys	AGA	Arg
	ATG	Met	ACG		AAG		AGG	
G	GTT	Val	GCT	Ala	GAT	Asp	GGT	Gly
	GTC		GCC		GAC		GGC	
	GTA		GCA		GAA	Glu	GGA	
	GTG		GCG		GAG		GGG	

Example:

ATGAAACAGTTGCTGGAAGCCGGTGTTCACCTTC
 TTATCAGAACTTTAGATGCCAGCGCTCACATA
 TTAGAAGATATGATTCAAAGTGGAAATGCATTTT
 CTCGAACAAATGCTAGATGTAGGTGTTCAATTTT
 TTACAGTCAATGCTTGAAGCTGGTGTTCACCTTT
 TTAGAAGCACTTTTAGAGGCTGGTGTTCATTTT
 TTGGAAGAAATGATGGAAGCGGGTATTCATTTT
 TTGGAACAAATGATGGAAGCAGGAGTCCATTTT
 ATTCAGCAATTATTGGAAGCAGGAGTCCATCTG
 ATTGAAGAATTGTTAGAAGCTGGCGTGCATTTT
 TTAGAACAAATGTTAGATGCAGGTGTACATTTT
 TTACA AAAAATGATTGAAGCTGGTGTTCATTTT
 CTATCAGAAATGATGGAAGCTGGTGTTCATTTT
 CTGCCGCAAATGCTGGAAGCCGGTGTCCATTTT
 TTAGAAGAAATGATGGAAGCAGGGTCCATTTT
 TTAGCAGAAATTAAGAAAGCGGGCGTTCATTTT

Degenerating synonymous 3rd codon positions

	T		C		A		G	
T	TTY	Phe	TCN	Ser	TAY	Tyr	TGY	Cys
	TTY		TCN		TAY		TGY	
	TTN	Leu	TCN		TAR	Ter	TGR	Ter
	TTN		TCN		TAR		TGG	Trp
C	CTN	Leu	CCN	Pro	CAY	His	CGN	Arg
	CTN		CCN		CAY		CGN	
	CTN		CCN		CAR	Gln	CGN	
	CTN		CCN		CAR		CGN	
A	ATH	Ile	ACN	Thr	AAV	Asn	AGN	Ser
	ATH		ACN		AAV		AGN	
	ATH		ACN		AAR	Lys	AGN	
	ATG	Met	ACN		AAR		AGN	Arg
G	GTN	Val	GCN	Ala	GAY	Asp	GGN	Gly
	GTN		GCN		GAY		GGN	
	GTN		GCN		GAR	Glu	GGN	
	GTN		GCN		GAR		GGN	

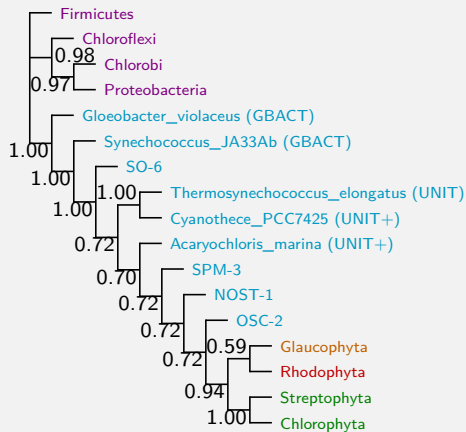
Example:

ATGAARCARTTNCTNGARGCNGGNGTNCAYTTY
 TTNTCNGARCTNTTNGAYGCNAGNGCNCAYATH
 TTNGARGAYATGATHCARAGNNGNATGCAITTY
 CTNGARCARATGCTNGAYGTNGGNGTNCAYTTY
 TTNCARTCNATGCTNGARGCNGGNGTNCAYTTY
 TTNGARGCNCTNTTNGARGCNGGNGTNCAYTTY
 TTNGARGARATGATGGARGCNGGNATHCAITTY
 TTNGARCARATGATGGARGCNGGNGTNCAYTTY
 ATHCARCARTTNTTNGARGCNGGNGTNCAYCTN
 ATHGARGARTTNTTNGARGCNGGNGTNCAYTTY
 TTNGARCARATGTTNGAYGCNNGGNGTNCAYTTY
 TTNCARAARATGATHGARGCNGGNGTNCAYTTY
 CTNTCNGARATGATGGARGCNGGNGCNCAYTTY
 CTNCCNCARATGCTNGARGCNGGNGTNCAYTTY
 TTNGARGARATGATGGARGCNGGNGTNCAYTTY
 TTNGCNGARTTNTCTNGARGCNGGNGTNCAYTTY

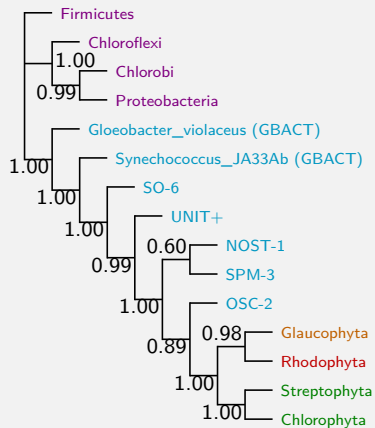
Degenerating synonymous 3rd codon positions

degenerate at 3rd pos.

(27.35% recoded)



cg75

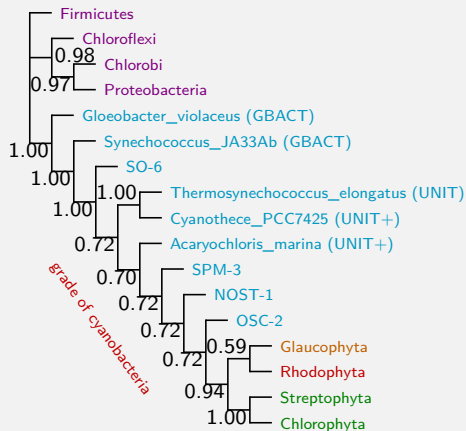


cg75_degenerate3

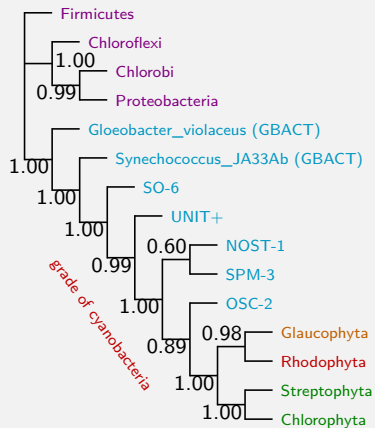
Degenerating synonymous 3rd codon positions

degenerate at 3rd pos.

(27.35% recoded)



cg75



cg75_degenerate3

Degenerating synonymous 3rd codon positions

- ▶ No "core" cyanobacteria monophyly.

Degenerating synonymous 3rd codon positions

- ▶ No "core" cyanobacteria monophyly.
→ Synonymous substitutions at 3rd codon position are not the only cause of discrepancy between cp75 and cg75.

Degenerating synonymous 3rd codon positions

- ▶ No "core" cyanobacteria monophyly.
→ Synonymous substitutions at 3rd codon position are not the only cause of discrepancy between cp75 and cg75.
- ▶ But there are synonymous substitutions also at 1st and 2nd position:

Degenerating synonymous 3rd codon positions

- ▶ No "core" cyanobacteria monophyly.
→ Synonymous substitutions at 3rd codon position are not the only cause of discrepancy between cp75 and cg75.
- ▶ But there are synonymous substitutions also at 1st and 2nd position:
e.g for Leu: (TTA ↔ CTA) → YTA
e.g for Arg: (CGA ↔ AGA) → MGA
e.g for Ser: (TCT ↔ AGA) → WST

Degenerating synonymous 1st and 2nd codon positions

	T		C		A		G	
T	TTT	Phe	TCT	Ser	TAT	Tyr	TGT	Cys
	TTC		TCC		TAC		TGC	
	TTA	Leu	TCA		TAA	Ter	TGA	Ter
	TTG		TCG		TAG		TGG	Trp
C	CTT	Leu	CCT	Pro	CAT	His	CGT	Arg
	CTC		CCC		CAC		CGC	
	CTA		CCA		CAA	Gln	CGA	
	CTG		CCG		CAG		CGG	
A	ATT	Ile	ACT	Thr	AAT	Asn	AGT	Ser
	ATC		ACC		AAC		AGC	
	ATA		ACA		AAA	Lys	AGA	
	ATG	Met	ACG		AAG		AGG	Arg
G	GTT	Val	GCT	Ala	GAT	Asp	GGT	Gly
	GTC		GCC		GAC		GGC	
	GTA		GCA		GAA	Glu	GGA	
	GTG		GCG		GAG		GGG	

Example:

```

ATGAAACAGTTGCTGGAAGCCGGTGTTCACTTC
TTATCAGAACTTTTAGATGCCAGCGCTCACATA
TTAGAAGATATGATTCAAAGTGGAATGCATTTT
CTCGAACAAATGCTAGATGTAGGTGTTCAATTTT
TTACAGTCAATGCTTGAAGCTGGTGTTCACTTT
TTAGAAGCACTTTTAGAGGCTGGTGTTCATTTT
TTGGAAGAAATGATGGAAGCGGGTATTCATTTT
TTGGAACAAATGATGGAAGCAGGAGTCCATTTT
ATTCAGCAATTATTGGAAGCAGGAGTCCATCTG
ATTGAAGAATTGTTAGAAGCTGGCGTGCATTTT
TTAGAACAAATGTTAGATGCAGGTGTACATTTT
TTACAAAAAATGATTGAAGCTGGTGTTCATTTT
CTATCAGAAATGATGGAAGCTGGTGTTCATTTT
CTGCCGCAAATGCTGGAAGCCGGTGTCCATTTT
TTAGAAGAAATGATGGAAGCAGGGGTCCATTTT
TTAGCAGAATTACTAGAAGCGGGCGTTCAATTTT

```


Degenerating synonymous 1st and 2nd codon positions

	T		C		A		G	
T	TTT	Phe	TCT	Ser	TAT	Tyr	TGT	Cys
	TTC		TCC		TAC		TGC	
	TTA	Leu	TCA		TAA	Ter	TGA	Ter
	TTG		TCCG		TAG		TGG	Trp
C	CTT	Leu	CCT	Pro	CAT	His	CGT	Arg
	CTC		CCC		CAC		CGC	
	CTA		CCA		CAA	Gln	CGA	
	CTG		CCG		CAG		CGG	
A	ATT	Ile	ACT	Thr	AAT	Asn	AGT	Ser
	ATC		ACC		AAC		AGC	
	ATA		ACA		AAA	Lys	AGA	
	ATG	Met	ACG		AAG		AGG	Arg
G	GTT	Val	GCT	Ala	GAT	Asp	GGT	Gly
	GTC		GCC		GAC		GGC	
	GTA		GCA		GAA	Glu	GGA	
	GTG		GCG		GAG		GGG	

Example:

ATGAAACAGTTGCTGGAAGCCGGTGTTCACTTC
 TTATCAGAACTTTAGATGCCAGCGCTCACATA
 TTAGAAGATATGATTCAAAGTGGAAATGCATTTT
 CTCGAACAAATGCTAGATGTAGGTGTTCAATTTT
 TTACAGTCAATGCTTGAAGCTGGTGTTCACTTT
 TTAGAAGCACTTTTAGAGGCTGGTGTTCATTTT
 TTGGAAGAAATGATGGAAGCGGGTATTCATTTT
 TTGGAACAAATGATGGAAGCAGGAGTCCATTTT
 ATTCAGCAATTATTGGAAGCAGGAGTCCATCTG
 ATTGAAGAAATTGTTAGAAGCTGGCGTGCATTTT
 TTAGAACAAATGTTAGATGCAGGTGTACATTTT
 TTACAAAAAATGATTGAAGCTGGTGTTCATTTT
 CTATCAGAAATGATGGAAGCTGGTGTTCATTTT
 CTGCCGCAAATGCTGGAAGCCGGTGTCCATTTT
 TTAGAAGAAATGATGGAAGCAGGGTCCATTTT
 TTAGCAGAACTACTAGAAGCGGGCGTTCAATTTT

Degenerating synonymous 1st and 2nd codon positions

	T		C		A		G	
T	TTT	Phe	WST	Ser	TAT	Tyr	TGT	Cys
	TTC		WSC		TAC		TGC	
	YTA	Leu	WSA		TAA	Ter	TGA	Ter
	YTG		WSG		TAG		TGG	Trp
C	YTT	Leu	CCT	Pro	CAT	His	MGT	Arg
	YTC		CCC		CAC		MGC	
	YTA		CCA		CAA	Gln	MGA	
	YTG		CCG		CAG		MGG	
A	ATT	Ile	ACT	Thr	AAT	Asn	WST	Ser
	ATC		ACC		AAC		WSC	
	ATA		ACA		AAA	Lys	MGA	
	ATG	Met	ACG		AAG		MGG	
G	GTT	Val	GCT	Ala	GAT	Asp	GGT	Gly
	GTC		GCC		GAC		GGC	
	GTA		GCA		GAA	Glu	GGA	
	GTG		GCG		GAG		GGG	

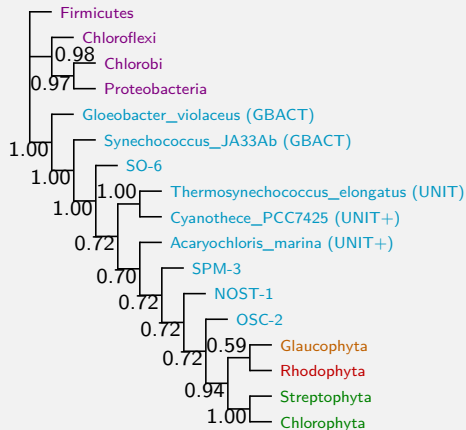
Example:

ATGAAACAGY**Y**TG**Y**TGGAAGCCGGTGTTC**A**CTTC
 YTA**W**SAGAA**Y**TT**Y**TAGATGCC**W**SCGCTCACATA
 YTAGAAGATATGATTCAA**W**STGGAATGCATTTT
 YTCGAACAAAT**G**YTAGATGTAGGTGTTCATTTT
 YTACAG**W**SAAT**G**YTTGAAGCTGGTGTTC**A**CTTT
 YTAGAAGCA**Y**TT**Y**TAGAGGCTGGTGTTCATTTT
 YTGGAAGAAATGATGGAAGCGGGTATTCATTTT
 YTGGAACAAATGATGGAAGCAGGAGTCCATTTT
 ATTCAGCAA**Y**T**A**YTGGAAGCAGGAGTCCAT**Y**TG
 ATTGAAGAA**Y**T**G**YTAGAAGCTGGCGTGCATTTT
 YTAGAACAAAT**G**YTAGATGCAGGTGTACATTTT
 YTACAAAAAATGATTGAAGCTGGTGTTCATTTT
 YTA**W**SAGAAAATGATGGAAGCTGGTGTTCATTTT
 YTGCCGCAAAT**G**YTGGAAGCCGGTGTCCATTTT
 YTAGAAGAAATGATGGAAGCAGGGGTCCATTTT
 YTAGCAGAA**Y**T**A**YTAGAAGCGGGCGTTCATTTT

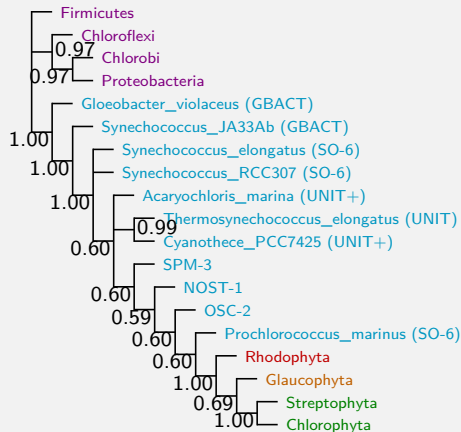
Degenerating synonymous 1st and 2nd codon positions

degenerate at 1st and 2nd pos.

(7.62% recoded)

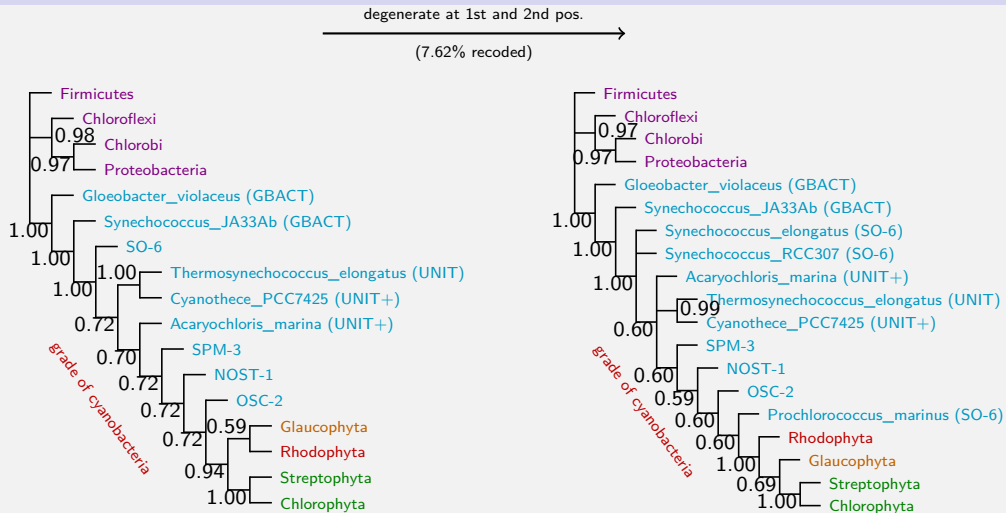


cg75



cg75_degenerate12

Degenerating synonymous 1st and 2nd codon positions



Degenerating synonymous 1st and 2nd codon positions

- ▶ Again no "core" cyanobacteria monophyly.

Degenerating synonymous 1st and 2nd codon positions

- ▶ Again no "core" cyanobacteria monophyly.
→ Synonymous substitutions at 1st and 2nd codon position are not the only cause of discrepancy between cp75 and cg75.

Degenerating synonymous 1st and 2nd codon positions

- ▶ Again no "core" cyanobacteria monophyly.
→ Synonymous substitutions at 1st and 2nd codon position are not the only cause of discrepancy between cp75 and cg75.
- ▶ Let's try to neutralize all synonymous substitutions. . .

Degenerating all synonymous codon positions

	T		C		A		G	
T	TTT	Phe	TCT	Ser	TAT	Tyr	TGT	Cys
	TTC		TCC		TAC		TGC	
	TTA	Leu	TCA		TAA	Ter	TGA	Ter
	TTG		TCG		TAG		TGG	Trp
C	CTT	Leu	CCT	Pro	CAT	His	CGT	Arg
	CTC		CCC		CAC		CGC	
	CTA		CCA		CAA	Gln	CGA	
	CTG		CCG		CAG		CGG	
A	ATT	Ile	ACT	Thr	AAT	Asn	AGT	Ser
	ATC		ACC		AAC		AGC	
	ATA		ACA		AAA	Lys	AGA	
	ATG	Met	ACG		AAG		AGG	Arg
G	GTT	Val	GCT	Ala	GAT	Asp	GGT	Gly
	GTC		GCC		GAC		GGC	
	GTA		GCA		GAA	Glu	GGA	
	GTG		GCG		GAG		GGG	

Example:

```

ATGAAACAGTTGCTGGAAGCCGGTGTTCACTTC
TTATCAGAACTTTTAGATGCCAGCGCTCACATA
TTAGAAGATATGATTCAAAGTGGAATGCATTTT
CTCGAACAAATGCTAGATGTAGGTGTTCAATTTT
TTACAGTCAATGCTTGAAGCTGGTGTTCACTTT
TTAGAAGCACTTTTAGAGGCTGGTGTTCATTTT
TTGGAAGAAATGATGGAAGCGGGTATTCATTTT
TTGGAACAAATGATGGAAGCAGGAGTCCATTTT
ATTCAGCAATTATTGGAAGCAGGAGTCCATCTG
ATTGAAGAATTGTTAGAAGCTGGCGTGCATTTT
TTAGAACAAATGTTAGATGCAGGTGTACATTTT
TTACAAAAAATGATTGAAGCTGGTGTTCATTTT
CTATCAGAAATGATGGAAGCTGGTGTTCATTTT
CTGCCGCAAATGCTGGAAGCCGGTGTCCATTTT
TTAGAAGAAATGATGGAAGCAGGGGTCCATTTT
TTAGCAGAATTACTAGAAGCGGGCGTTCAATTTT

```


Degenerating all synonymous codon positions

	T		C		A		G	
T	TTT	Phe	TCT	Ser	TAT	Tyr	TGT	Cys
	TTC		TCC		TAC		TGC	
	TTA	Leu	TCA		TAA	Ter	TGA	Ter
	TTG		TCCG		TAG		TGG	Trp
C	CTT	Leu	CCT	Pro	CAT	His	CGT	Arg
	CTC		CCC		CAC		CGC	
	CTA		CCA		CAA	Gln	CGA	
	CTG		CCG		CAG		CGG	
A	ATT	Ile	ACT	Thr	AAT	Asn	AGT	Ser
	ATC		ACC		AAC		AGC	
	ATA		ACA		AAA	Lys	AGA	
	ATG	Met	ACG		AAG		AGG	Arg
G	GTT	Val	GCT	Ala	GAT	Asp	GGT	Gly
	GTC		GCC		GAC		GGC	
	GTA		GCA		GAA	Glu	GGA	
	GTG		GCG		GAG		GGG	

Example:

ATGAAACAGTTGCTGGAAGCCGGTGTTCACCTTC
 TTATCAGAACTTTAGATGCCAGCGCTCACATA
 TTAGAAGATATGATTCAAAGTGGAAATGCATTTT
 CTCGAACAAATGCTAGATGTAGGTGTTCAATTTT
 TTACAGTCAATGCTTGAAGCTGGTGTTCACCTTT
 TTAGAAGCACTTTTAGAGGCTGGTGTTCATTTT
 TTGGAAGAAATGATGGAAGCGGGTATTCATTTT
 TTGGAACAAATGATGGAAGCAGGAGTCCATTTT
 ATTCAGCAATTATTGGAAGCAGGAGTCCATCTG
 ATTGAAGAATTGTTAGAAGCTGGCGTGCATTTT
 TTAGAACAAATGTTAGATGCAGGTGTACATTTT
 TTACAATAAATGATTGAAGCTGGTGTTCATTTT
 CTATCAGAAATGATGGAAGCTGGTGTTCATTTT
 CTGCCGCAATGCTGGAAGCCGGTGTCCATTTT
 TTAGAAGAAATGATGGAAGCAGGGTCCATTTT
 TTAGCAGAATTACTAGAAGCGGGCGTTCATTTT

Degenerating all synonymous codon positions

	T		C		A		G	
T	TTY	Phe	WSN	Ser	TAY	Tyr	TGY	Cys
	TTY		WSN		TAY		TGY	
	YTN	Leu	WSN		TAR	Ter	TGR	Ter
	YTN		WSN		TAR		TGG	Trp
C	YTN	Leu	CCN	Pro	CAY	His	MGN	Arg
	YTN		CCN		CAY		MGN	
	YTN		CCN		CAR	Gln	MGN	
	YTN		CCN		CAR		MGN	
A	ATH	Ile	ACN	Thr	AAV	Asn	WSN	Ser
	ATH		ACN		AAV		WSN	
	ATH		ACN		AAR	Lys	MGN	
	ATG	Met	ACN		AAR		MGN	Arg
G	GTN	Val	GCN	Ala	GAY	Asp	GGN	Gly
	GTN		GCN		GAY		GGN	
	GTN		GCN		GAR	Glu	GGN	
	GTN		GCN		GAR		GGN	

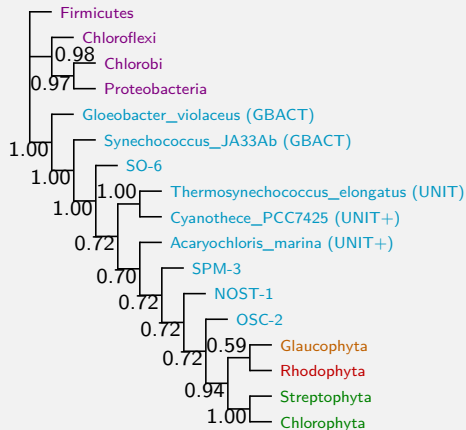
Example:

ATGAAR**C**AR**Y**T**N**Y**T**NGAR**G**C**N**GG**N**GT**N**CAY**T**TY
 Y**T**N**W**S**N**GAR**Y**T**N**Y**T**NGAY**G**C**N****W**S**N**G**C**N**C**AY**A**TH
 Y**T**NGAR**G**AYATGATH**C**AR**W**S**N**GG**N**ATG**C**AY**T**TY
 Y**T**NGAR**C**ARATGY**T**NGAY**G**T**N**GG**N**GT**N**CAY**T**TY
 Y**T**NCAR**W**S**N**ATGY**T**NGAR**G**C**N**GG**N**GT**N**CAY**T**TY
 Y**T**NGAR**G**C**N**Y**T**NY**T**NGAR**G**C**N**GG**N**GT**N**CAY**T**TY
 Y**T**NGAR**G**ARATGATGGAR**G**C**N**GG**N**ATH**C**AY**T**TY
 Y**T**NGAR**C**ARATGATGGAR**G**C**N**GG**N**GT**N**CAY**T**TY
 ATH**C**AR**C**AR**Y**T**N**Y**T**NGAR**G**C**N**GG**N**GT**N**CAY**Y**T**N**
 ATH**G**AR**G**AR**Y**T**N**Y**T**NGAR**G**C**N**GG**N**GT**N**CAY**T**TY
 Y**T**NGAR**C**ARATGY**T**NGAY**G**C**N**GG**N**GT**N**CAY**T**TY
 Y**T**NCAR**A**ARATGATH**G**AR**G**C**N**GG**N**GT**N**CAY**T**TY
 Y**T**N**W**S**N**GARATGATGGAR**G**C**N**GG**N**G**C**N**C**AY**T**TY
 Y**T**NC**C**N**C**ARATGY**T**NGAR**G**C**N**GG**N**GT**N**CAY**T**TY
 Y**T**NGAR**G**ARATGATGGAR**G**C**N**GG**N**GT**N**CAY**T**TY
 Y**T**NG**C**N**G**AR**Y**T**N**Y**T**NGAR**G**C**N**GG**N**GT**N**CAY**T**TY

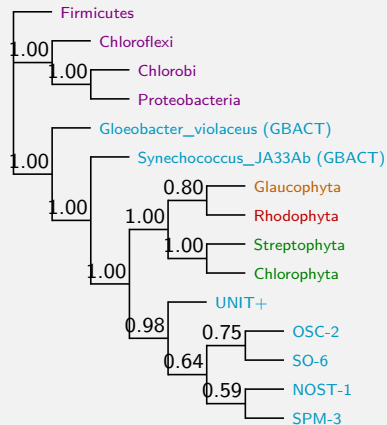
Degenerating all synonymous codon positions

degenerate at 1st, 2nd and 3rd pos.

(34.97% recoded)



cg75

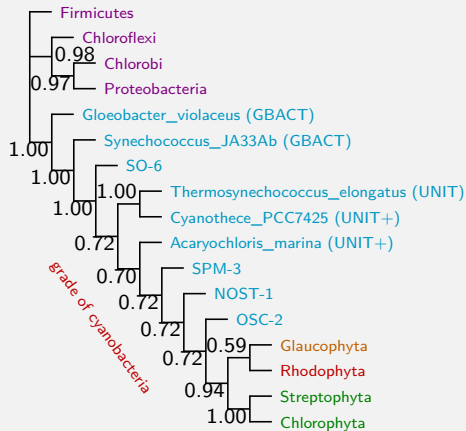


cg75_degenerate123

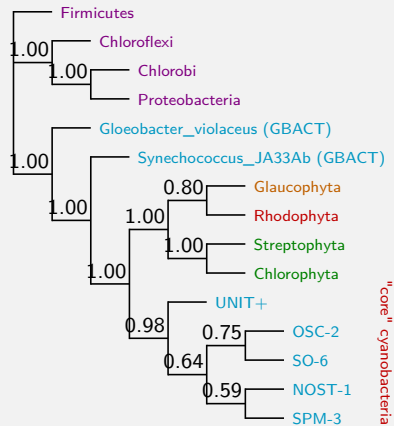
Degenerating all synonymous codon positions

degenerate at 1st, 2nd and 3rd pos.

(34.97% recoded)



cg75



cg75_degenerate123

Degenerating all synonymous codon positions

- ▶ Same topology as with the protein data

Degenerating all synonymous codon positions

- ▶ Same topology as with the protein data
- ▶ Synonymous substitutions are responsible for the incongruence between nucleotide and amino-acid data.

Degenerating all synonymous codon positions

- ▶ Same topology as with the protein data
- ▶ Synonymous substitutions are responsible for the incongruence between nucleotide and amino-acid data.
- ▶ This could be explained by effects of biases in codon usage and DNA composition:

Degenerating all synonymous codon positions

- ▶ Same topology as with the protein data
- ▶ Synonymous substitutions are responsible for the incongruence between nucleotide and amino-acid data.
- ▶ This could be explained by effects of biases in codon usage and DNA composition:
 - ▶ Codon preference / composition biases: "allowed" within a synonymy class (neutral at the protein level and above)

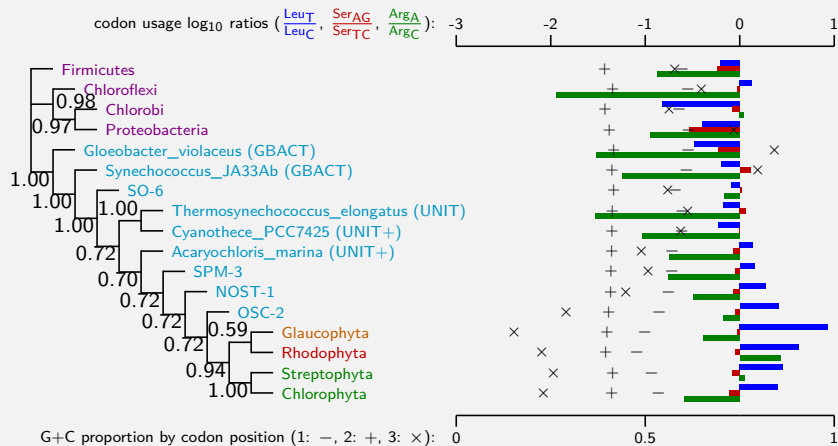
Degenerating all synonymous codon positions

- ▶ Same topology as with the protein data
- ▶ Synonymous substitutions are responsible for the incongruence between nucleotide and amino-acid data.
- ▶ This could be explained by effects of biases in codon usage and DNA composition:
 - ▶ Codon preference / composition biases: "allowed" within a synonymy class (neutral at the protein level and above)
 - Convergence between taxa sharing the same biases

Degenerating all synonymous codon positions

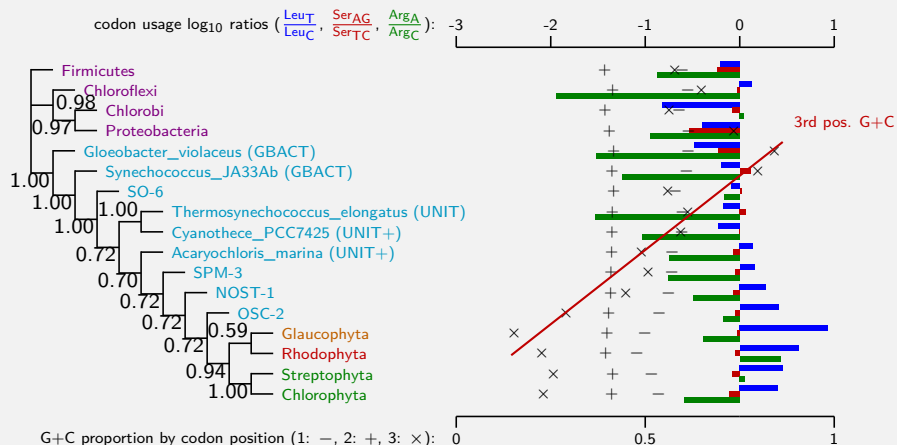
- ▶ Same topology as with the protein data
- ▶ Synonymous substitutions are responsible for the incongruence between nucleotide and amino-acid data.
- ▶ This could be explained by effects of biases in codon usage and DNA composition:
 - ▶ Codon preference / composition biases: "allowed" within a synonymy class (neutral at the protein level and above)
 - Convergence between taxa sharing the same biases
 - ▶ Let's visualize this. . .

Effects of composition and codon usage



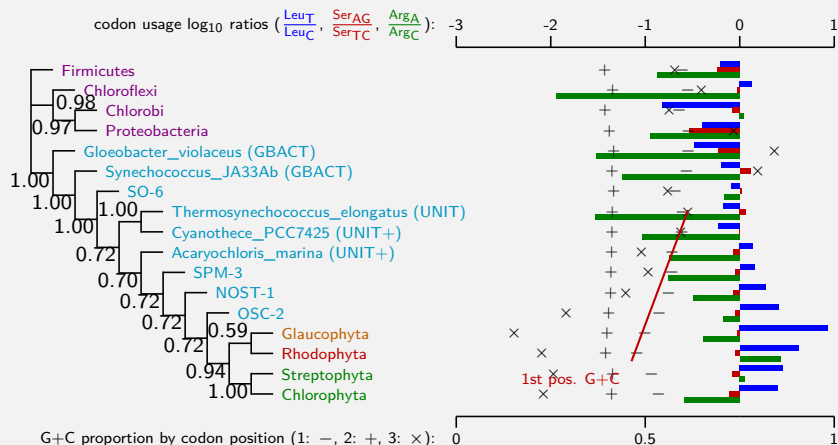
cg75

Effects of composition and codon usage



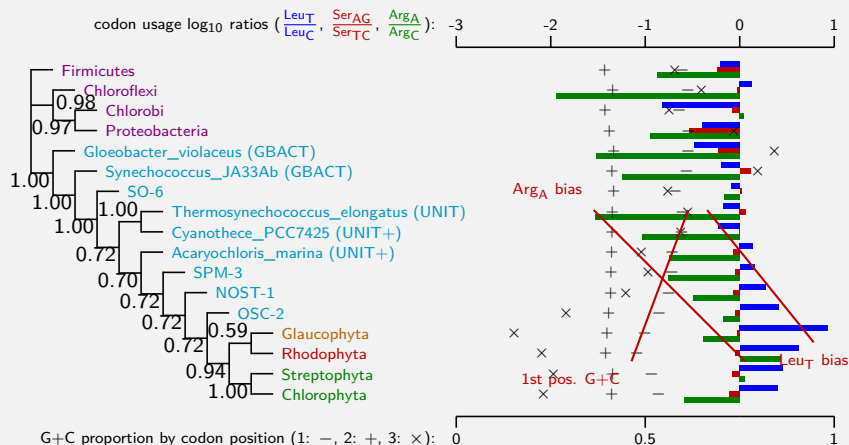
cg75

Effects of composition and codon usage



cg75

Effects of composition and codon usage



cg75

The conclusions

- ▶ Synonymous substitutions enable convergence in codon usage and composition biases.

The conclusions

- ▶ Synonymous substitutions enable convergence in codon usage and composition biases.
 - Nucleotide data are prone to reconstruction errors for large scale relationships (more time for convergence).

The conclusions

- ▶ Synonymous substitutions enable convergence in codon usage and composition biases.
 - Nucleotide data are prone to reconstruction errors for large scale relationships (more time for convergence).
 - We favour the hypothesis of a monophyletic "core" cyanobacteria clade sister to plastids (protein topology).

The conclusions

- ▶ Synonymous substitutions enable convergence in codon usage and composition biases.
 - Nucleotide data are prone to reconstruction errors for large scale relationships (more time for convergence).
 - We favour the hypothesis of a monophyletic "core" cyanobacteria clade sister to plastids (protein topology).
- ▶ But information potentially useful for small scale resolution is lost in translated / recoded data

The conclusions

- ▶ Synonymous substitutions enable convergence in codon usage and composition biases.
 - Nucleotide data are prone to reconstruction errors for large scale relationships (more time for convergence).
 - We favour the hypothesis of a monophyletic "core" cyanobacteria clade sister to plastids (protein topology).
- ▶ But information potentially useful for small scale resolution is lost in translated / recoded data
- ▶ The use of more advanced models could turn misleading signal into useful signal.

The conclusions

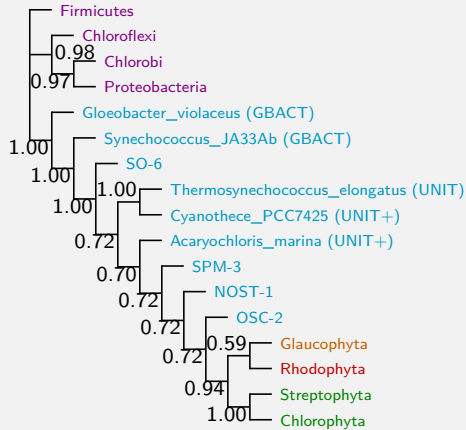
- ▶ Synonymous substitutions enable convergence in codon usage and composition biases.
 - Nucleotide data are prone to reconstruction errors for large scale relationships (more time for convergence).
 - We favour the hypothesis of a monophyletic "core" cyanobacteria clade sister to plastids (protein topology).
- ▶ But information potentially useful for small scale resolution is lost in translated / recoded data
- ▶ The use of more advanced models could turn misleading signal into useful signal.
 - Could we then get both large and small scale resolution using un-recoded nucleotide data?

Thanks for your attention

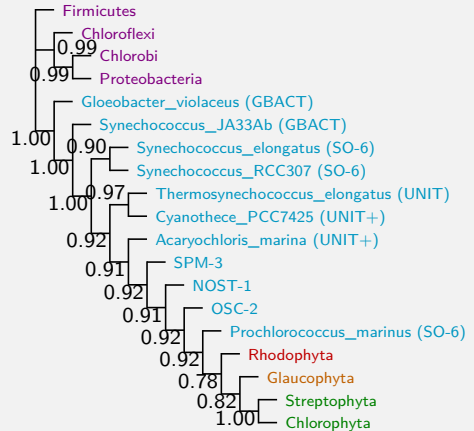
- ▶ This work was supported by a Fundação para a Ciência e a Tecnologia (FCT, Portugal) grant to Cymon J. Cox, Centro de Ciencias do Mar (CCMAR) - CIMAR-Lab. Assoc., (PTDC/BIA-BCM/099565/2008).
- ▶ Contact: blaise.li@normalesup.org

Effects of serine synonymous substitutions

degenerate serines (at pos. 1 and 2)

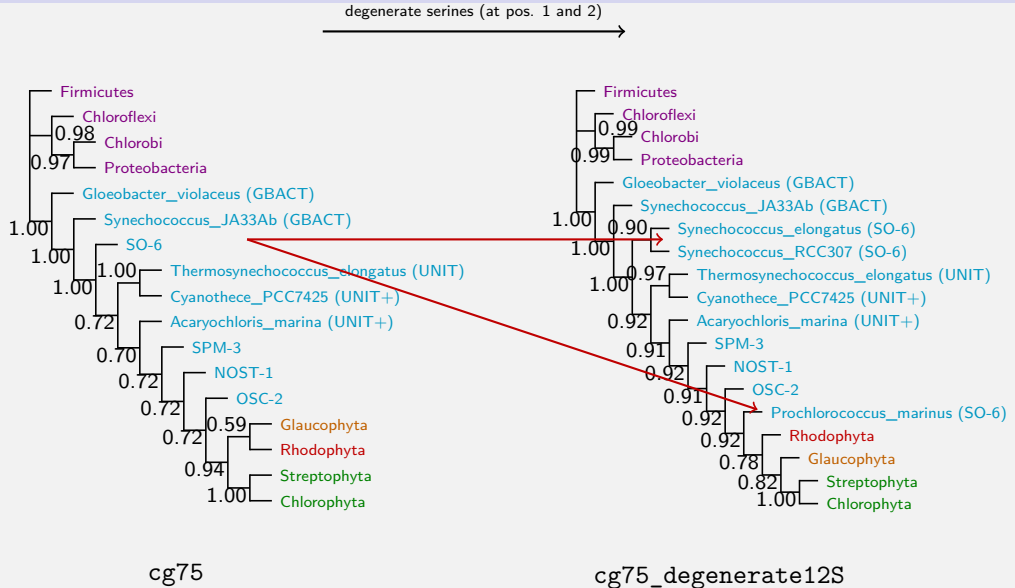


cg75



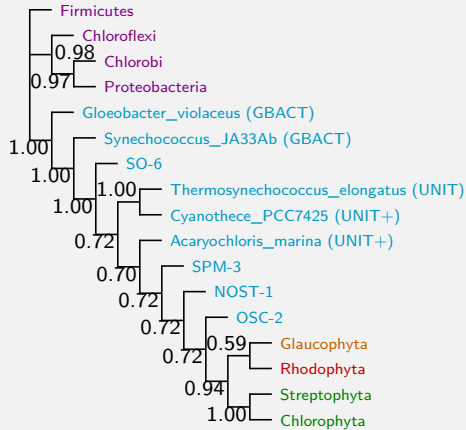
cg75_degenerate12S

Effects of serine synonymous substitutions

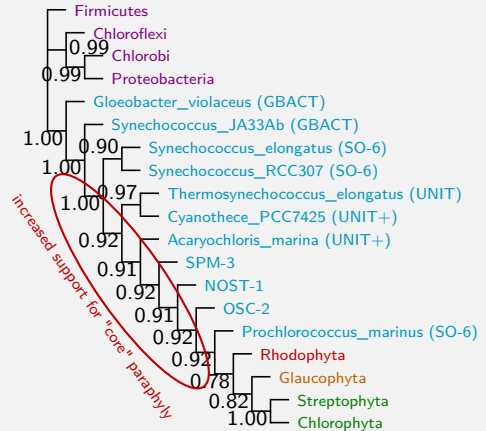


Effects of serine synonymous substitutions

degenerate serines (at pos. 1 and 2)



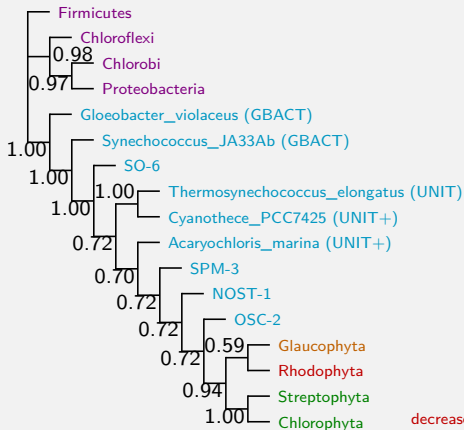
cg75



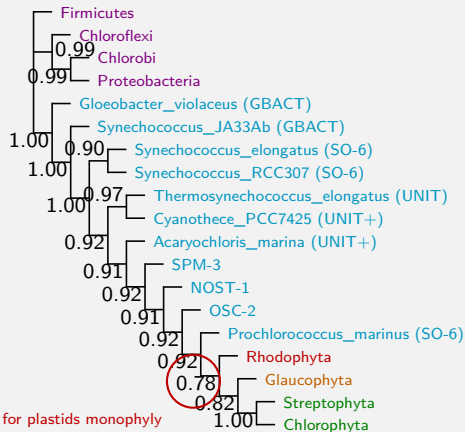
cg75_degenerate12S

Effects of serine synonymous substitutions

degenerate serines (at pos. 1 and 2)



cg75



cg75_degenerate12S

Effects of serine synonymous substitutions

- ▶ Serine signal:

Effects of serine synonymous substitutions

- ▶ Serine signal:
 - ▶ contributes to SO-6 and plastid monophyly

Effects of serine synonymous substitutions

- ▶ Serine signal:
 - ▶ contributes to SO-6 and plastid monophyly
 - ▶ decreases the support of "core" cyanobacteria paraphyly

Effects of serine synonymous substitutions

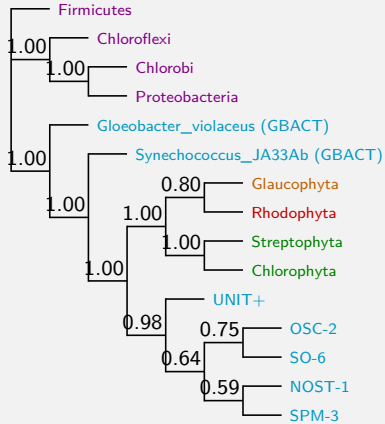
- ▶ Serine signal:
 - ▶ contributes to SO-6 and plastid monophyly
 - ▶ decreases the support of "core" cyanobacteria paraphyly
- ▶ Serine signal is not present in translated datasets

Effects of serine synonymous substitutions

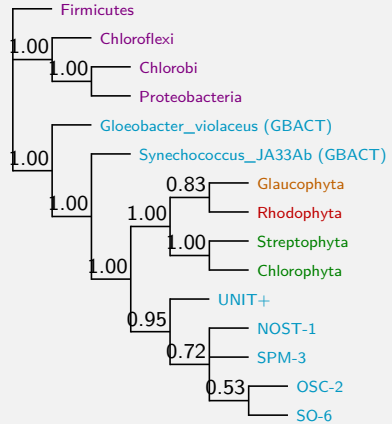
- ▶ Serine signal:
 - ▶ contributes to SO-6 and plastid monophyly
 - ▶ decreases the support of "core" cyanobacteria paraphyly
- ▶ Serine signal is not present in translated datasets
- ▶ Maybe we should keep this signal in the nucleotide dataset?

Effects of serine synonymous substitutions

restore serine codons →



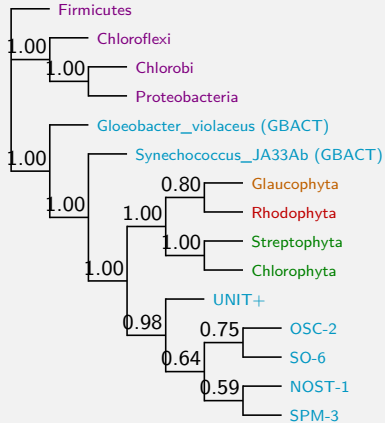
cg75_degenerateLRS3



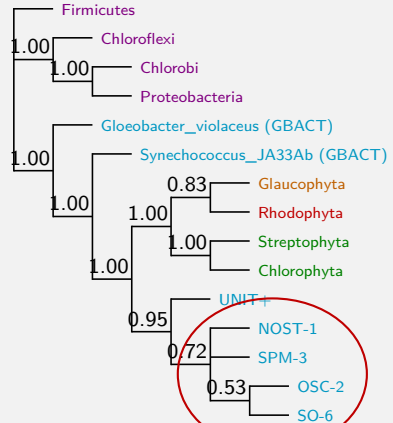
cg75_degenerateLR3

Effects of serine synonymous substitutions

restore serine codons →



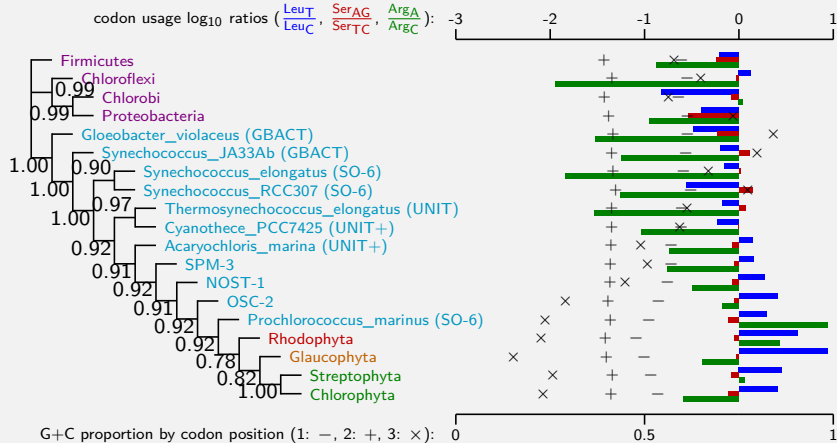
cg75_degenerateLRS3



unreliable relationships

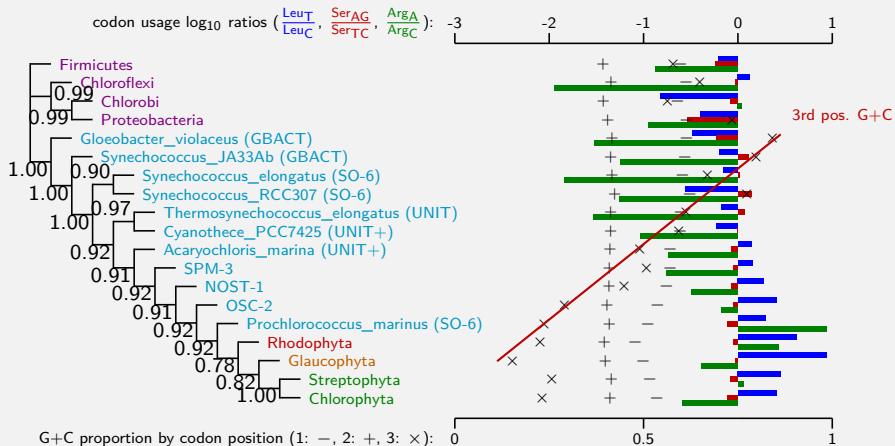
cg75_degenerateLR3

Effects of composition and codon usage



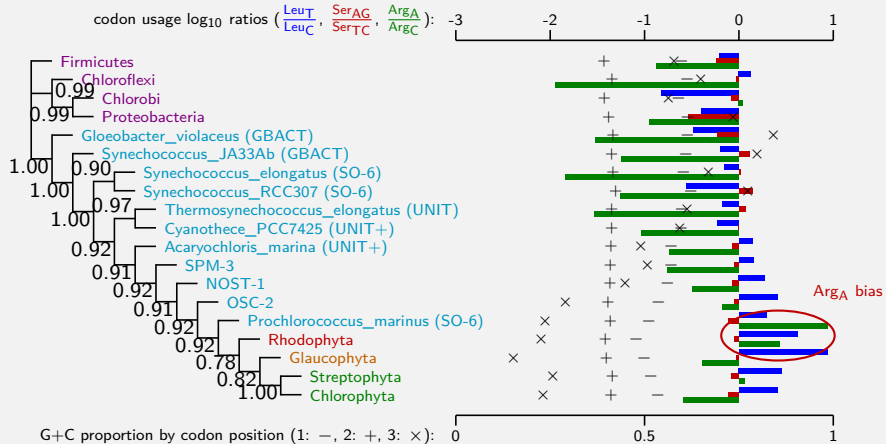
cg75_degenerate12S

Effects of composition and codon usage



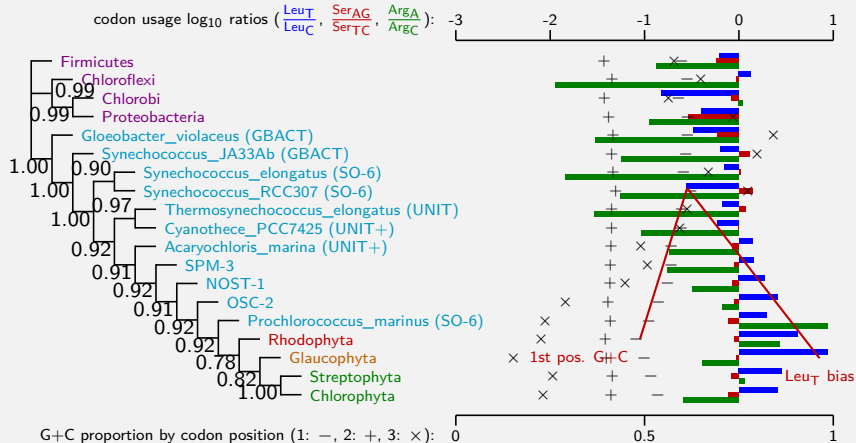
cg75_degenerate12S

Effects of composition and codon usage



cg75_degenerate12S

Effects of composition and codon usage



cg75_degenerate12S