**Supplementary Material for:**

**Multidomain ribosomal protein trees and the planctobacterial origin of neomura (eukaryotes, archaebacteria)**

Thomas Cavalier-Smith and Ema E-Yung Chao

**Contents:** Supplementary Figures S1-S17. pp. 1-23

## Figure S1. Site-heterogeneous PhyloBayes GTR-CAT (4 gamma rates) tree for 26 ribosomal proteins from 143 eukaryotes representing all the most divergent lineages.

Alignment of 4156 amino acids; the number included for each taxon is shown after the @. Support values are posterior probabilities. 105,136 trees were summed after removing the first 20% as burnin: maxdiff 0.29708. The tree is rooted between Eozoa and neokaryotes
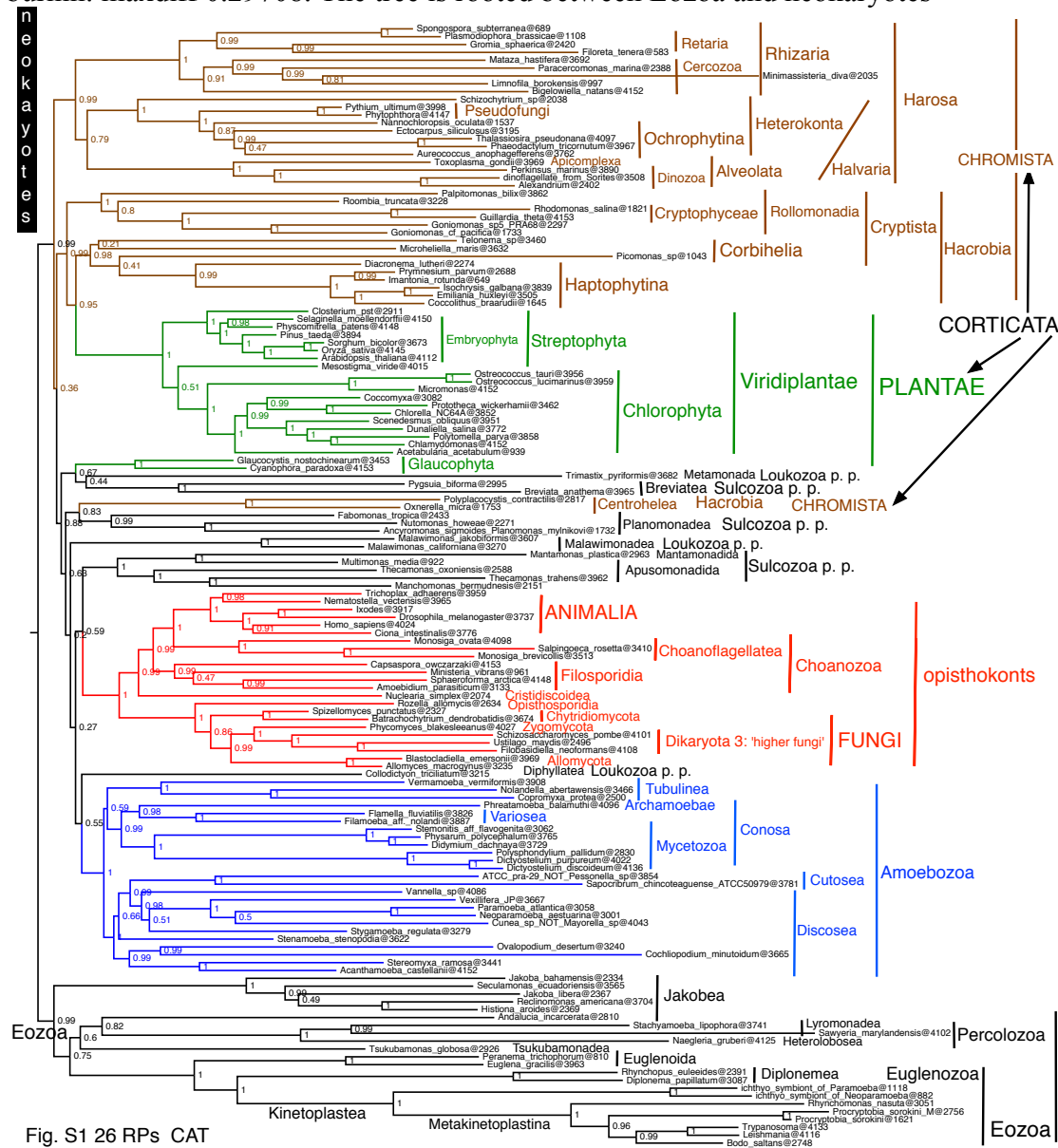


Fig. S1 26 RPs CAT

**Fig. S2. Site-homogeneous RAxML PROTGAMMALGF (4 gamma rates) tree for 26 ribosomal proteins from 151 eukaryotes representing all the most divergent lineages.**

Alignment of 4156 amino acids; the number included for each taxon is shown after the @. The tree is rooted between Eozoa and neokaryotes. Support values are percentages for 100 pseudoreplicates.
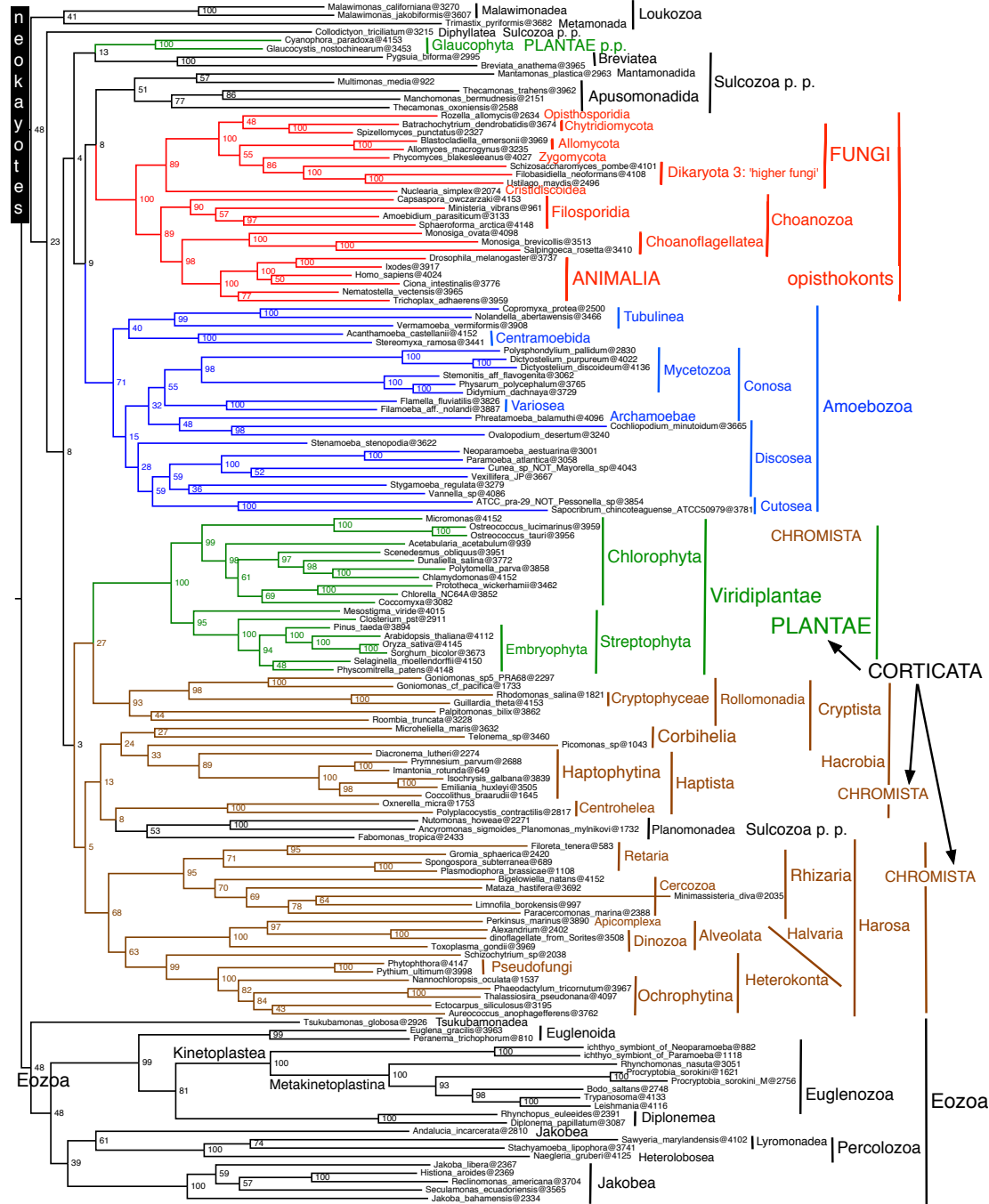


Fig. S2 26 RPs ML

**Figure S3. Site-heterogeneous PhyloBayes GTR-CAT (4 gamma rates) tree for 51 ribosomal proteins from 60 archaebacteria representing all the most divergent lineages.**
This is chain 2 which placed 'Nanohaloarchaea' in a different position from chain 1 (see Fig. 4). (49, 886 trees summed after removing 40% as burnin).
Support values are posterior probabilities.

**Figure S4. Site-heterogeneous PhyloBayes GTR-CAT (4 gamma rates) tree for 26 ribosomal proteins from 60 archaebacteria representing all the most divergent lineages.**
Support values are posterior probabilities. after removing 40% of the trees as burnin the remaining 55,580 trees from two chains were summed: maxdiff 0.0669948.



Fig. S4 26 RPs CAT

**Fig. S5. Site-homogeneous RAxML PROTGAMMALGF (4 gamma rates) tree for 26 ribosomal proteins from 60 archaebacteria representing all the most divergent lineages.**

Support values are percentages for 100 pseudoreplicates.



Fig. S5 26 RPs ML

**Figure S6. Site-homogeneous RAxML PROTGAMMALGF (4 gamma rates) tree for 26 ribosomal proteins from 151 eubacteria representing all the most divergent lineages.** Support values are percentages for 100 pseudoreplicates.
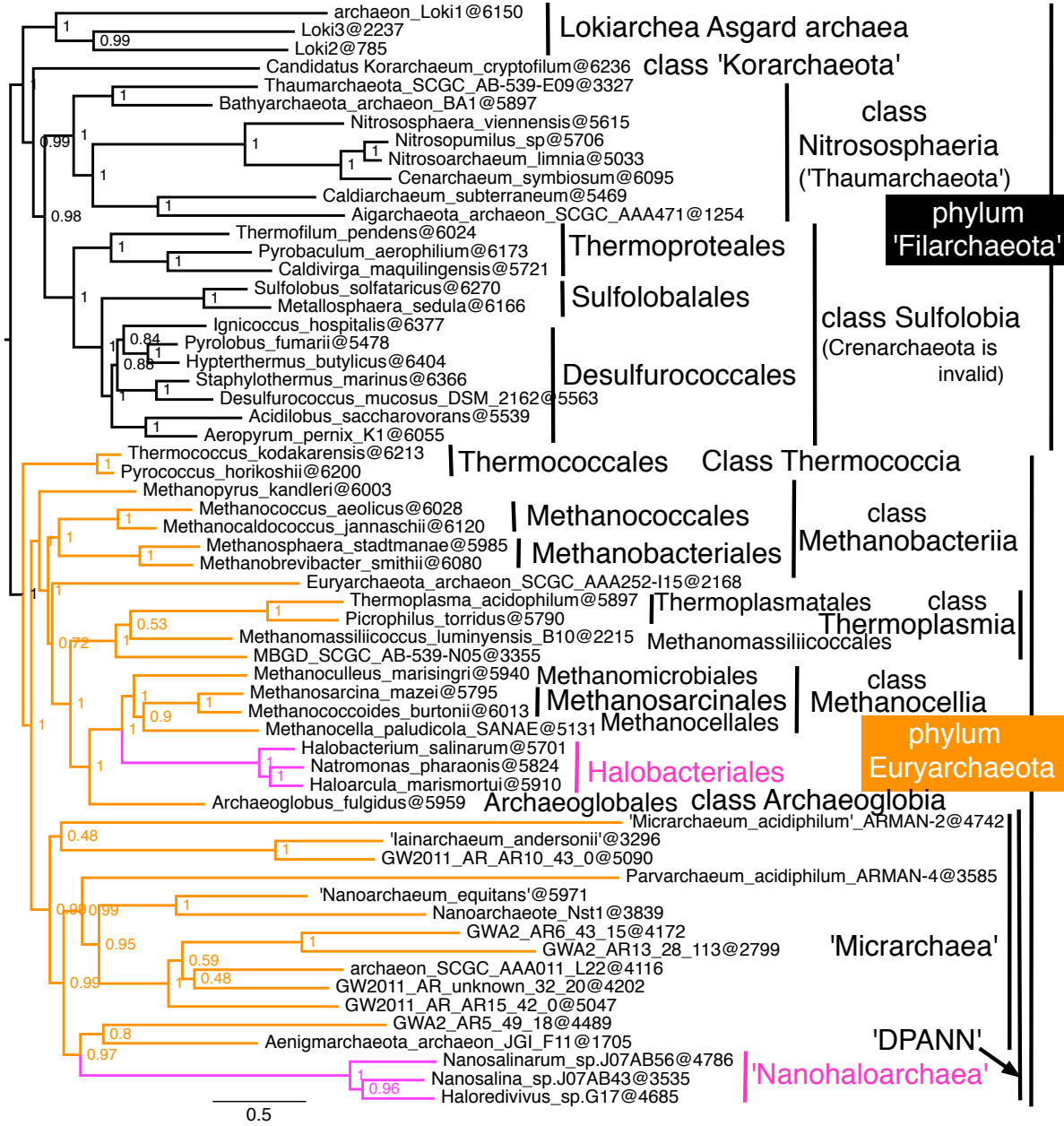
7

**Fig. S7. Site-heterogeneous PhyloBayes GTR-CAT (4 gamma rates) tree for 26 ribosomal proteins from 203 neomura representing all the most divergent lineages.**
Support values are posterior probabilities.

**Figure S8. Site-homogeneous RAxML PROTGAMMALGF (4 gamma rates) tree for 26 ribosomal proteins from 203 neomura representing all the most divergent lineages.** Support values are percentages for 100 pseudoreplicates.



**Fig. S9. Site-heterogeneous prokaryote PhyloBayes CAT-GTR tree for 51 ribosomal proteins from 60 archaebacteria and 26 ribosomal proteins from 151 eubacteria representing all the most divergent lineages, including chloroplasts.** Support values are posterior probabilities. See next page:

Fig. S9

**Figure S10. Site-homogeneous prokaryote RAxML PROTGAMMALGF (4 gamma rates) tree for 26 ribosomal proteins from 60 archaebacteria and 151 eubacteria representing all the most divergent lineages.**
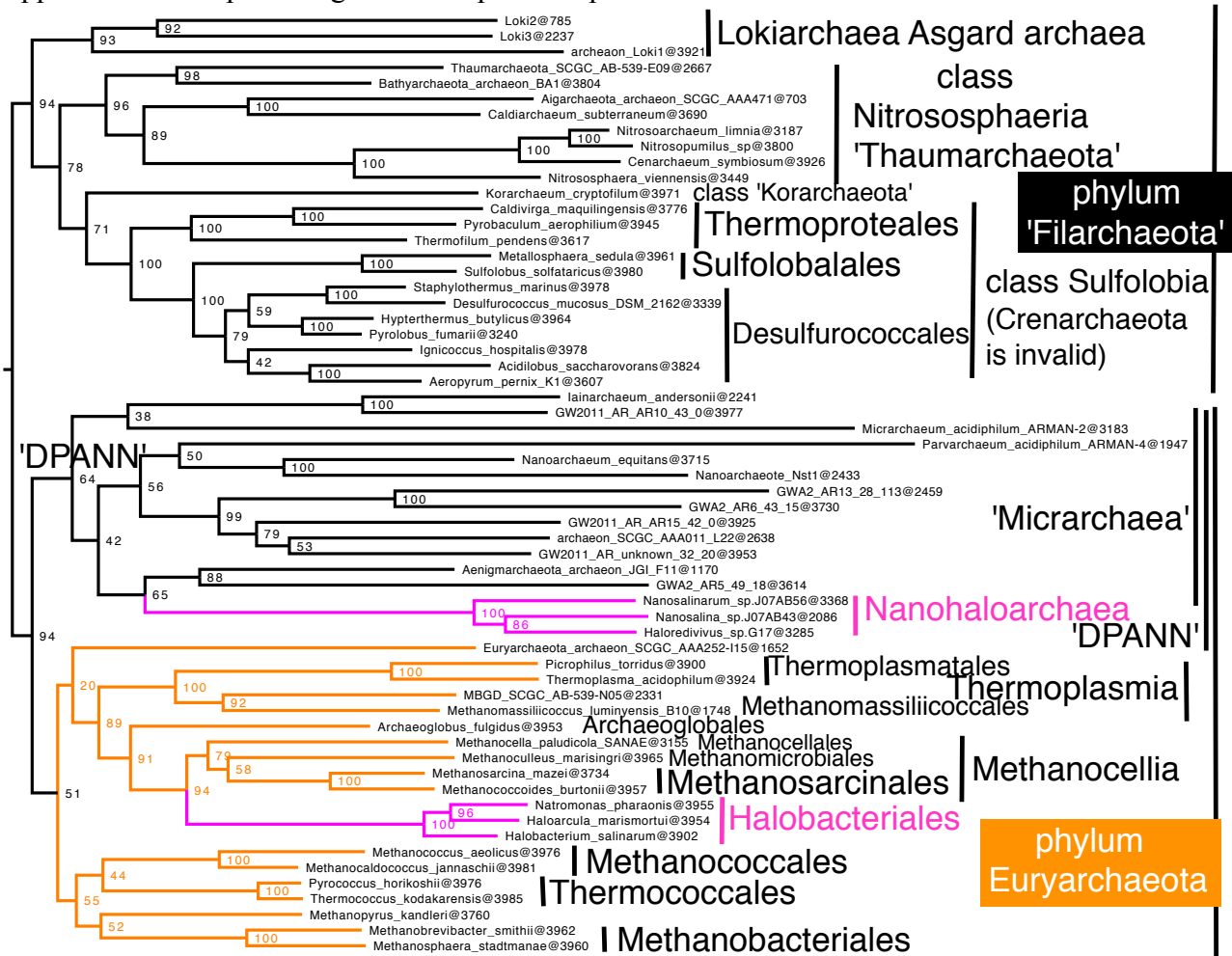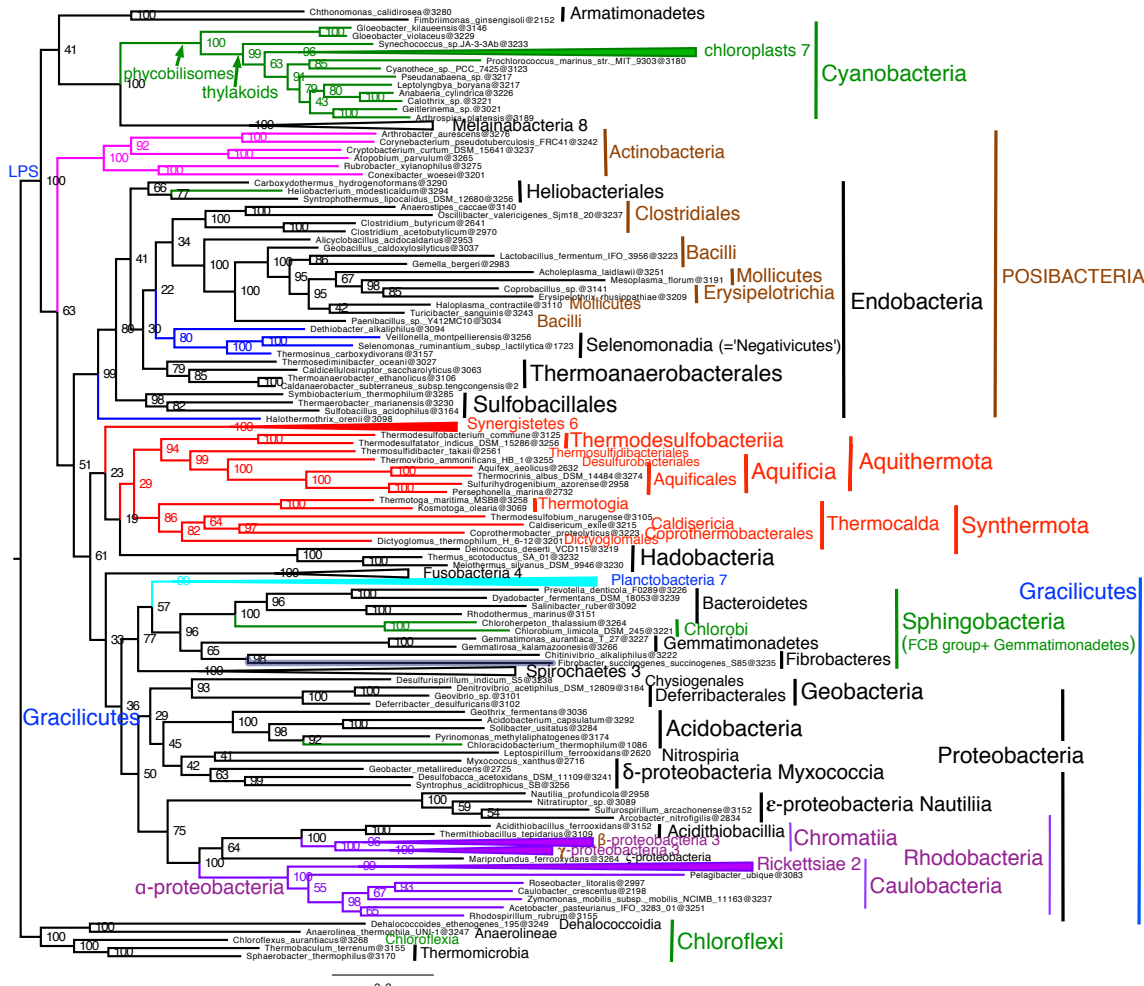
Support values are percentages for 100 pseudoreplicates. See next page:

Fig. S10

**Fig. S11. Site-heterogeneous 2-domain PhyloBayes CAT-GTR tree for 26 ribosomal proteins from 143 eukaryotes and 151 eubacteria representing all the most divergent lineages.**
Support values are posterior probabilities. Consensus tree for two chains: 38,969 trees were summed after removing 40% as burnin. The chains converged on the same topology except for a few nodes with .5 or less support: Max diff 1. See next page:

Chloroflexi

Armatimonadetes

Negativicutes

Endobacteria

Clostridia sensu stricto    POSIBACTERIA p. p.

Bacilli

+ Mollicutes

Thermocalda    Synthermota

Aquithermota

Synergistetes    Synthermota

Hadobacteria

EUBACTERIA

Actinobacteria    POSIBACTERIA p. p.

Fusobacteria

Melainabacteria

Cyanobacteria

chloroplasts

Cyanobacteria

Chysiogenales

Geobacteria p. p.

Deferribacterales    Nitrospira

δ-proteobacteria    Myxococcia

Acidobacteria

Proteobacteria

ε-proteobacteria    Geobacteria p. p.

Nautilia

Rhodobacteria

Gracilicutes

Spirochaetes

Sphingobacteria 10
(FCB group + Gemmatimonadetes)

Planctochlora

Planctobacteria

Eozoa

Percolozoa

Jakobea

EUKARYOTA

Euglenozoa

Neokaryota

Rollomonadia

Hacrobia p. p.

Corbihelia

Planomonadea    Sulcozoa p.p.

Corbihelia

Haptophytina

CHROMISTA

Retaria

Rhizaria

Cercozoa    Harosa

Heterokonta

Halvaria

Alveolata

Glaucophyta

PLANTAE

Streptophyta

Viridiplantae

Chlorophyta

Metamonada

Centrohelea

Diphyllatea

Malawimonadea

Mantamonadea

Apusomonadida

Sulcozoa p.p

Breviatea

Planomonadea

Choanoflagellatea

ANIMALIA

opisthokonts

Filosporidia    Choanozoa

Chytridiomycota

FUNGI

Zygomycota

Dikarya: higher fungi

Allomycota

Tubulinea

Archamoebae

Himatismenida

Variosea

Mycetozoa

Amoebozoa

Discosea

Cutosea

**Figure S12. Site-homogeneous 2-domain RAxML PROTGAMMALGF (4 gamma rates) tree for 26 ribosomal proteins from 143 eukaryotes and 151 eubacteria eubacteria representing all the most divergent lineages.**
Support values are percentages for 100 pseudoreplicates.

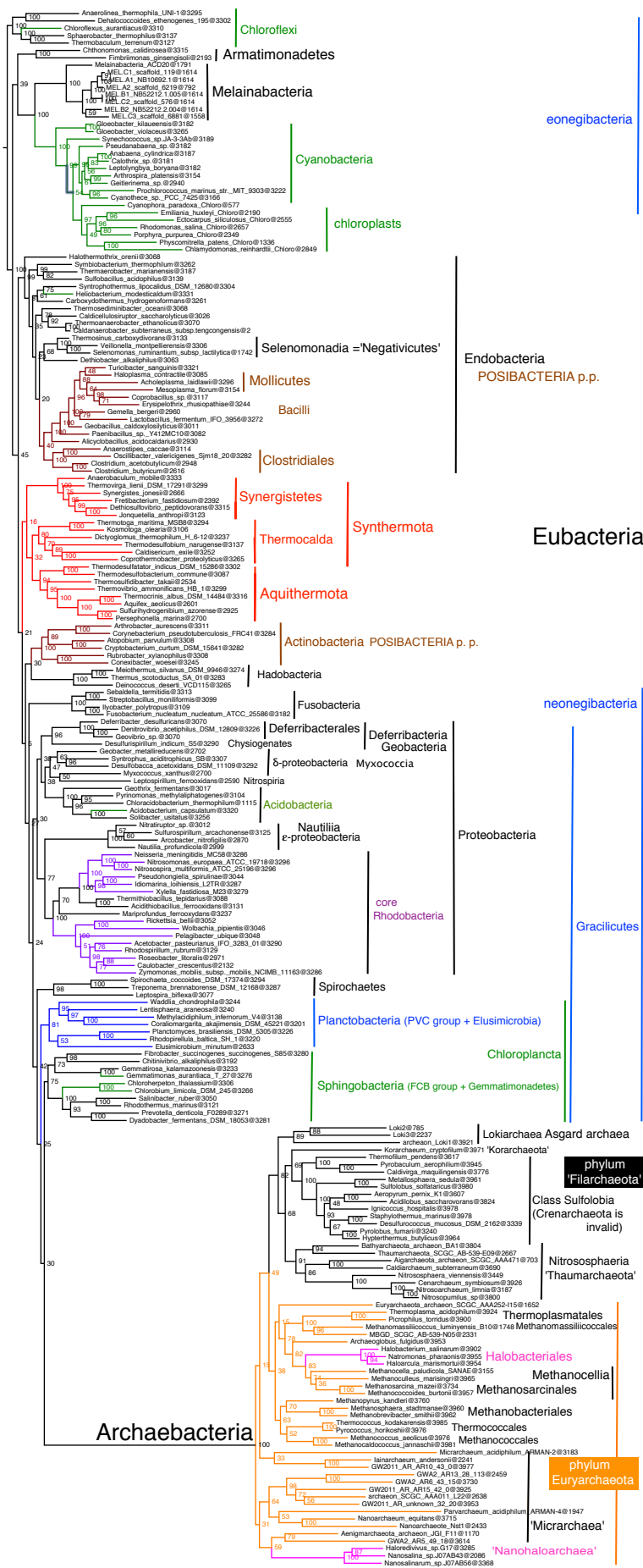**Fig. S13. Site-heterogeneous universal 3-domain PhyloBayes CAT-GTR tree for 26 ribosomal proteins from 143 eukaryotes, 60 archaebacteria, and 151 eubacteria representing all the most divergent lineages.**

Support values are posterior probabilities. Two chains were summed (28,777 trees) after removing the first 40% pre-log-likelihood plateau trees as burnin; despite clear plateauing a few deep-branching topological differences (PP 0.5 or less only) remained between the two chains, making maxdiff 1. This consensus tree and chain 1 are the only trees that wrongly placed Thermodesulfobacteriaceae inside Proteobacteria, implying that in chain 1 PhyloBayes got stuck in the wrong topology and did not cope well with the large number of taxa. Chain 2 had maximum support for their being sisters of Aquificia, their usual position; the corresponding ML tree (Fig. S14) confirms this with 89% support. See next page:

Chloroflexi

Armatimonadetes

Melainabacteria

Cyanobacteria

chloroplasts

Actinobacteria

POSIBACTERIA

Halanaerobiia

Selenomonadia

Endobacteria    mostly monoderm POSIBACTERIA

Clostridia

Bacilli

POSIBACTERIA

Mollicutes

Mollicutes

Fusobacteria

Spirochaetes

Planctobacteria

planctochlora

Sphingobacteria

Deferribacteria  Geobacteria

Deferribacterales  Chrysiogenales

Thermodesulfobacteriia Aquithermota p. p.

δ-proteobacteria Myxococcia

Proteobacteria

Acidobacteria

Gracilicutes
(=hydrobacteria)

ε-proteobacteria Nautiliia

Proteobacteria

Rhodobacteria

EUBACTERIA

Synergistetes  Synthermota

Hadobacteria

Thermotoga

Dictyoglomia/Caldisericia   Thermocalda  Synthermota

Aquificia Aquithermota p. p.

ARCHAEBACTERIA

'DPANN'

'Nanohaloarchaea'

phylum
Euryarchaeota

Thermococcales

Methanococcales

Methanobacteriia

Methanobacterales

Thermoplasmatales

Methanomicrobiales

Methanocellia

Methanosarcinales

Halobacteriales

Archaeoglobales

Nitrososphaeria
'Thaumarchaeota'

phylum 'Filarchaeota'

Class Sulfolobia
Crenarchaeota ia is invalid

'Korarchaeota'

Lokiarchaea Asgard archaea

inflated neomuran stem

Percolozoa

Discicristata

Eozoa

inflated eukaryote stem

Euglenozoa

Jakobea

Rollomonadia

Cryptista

Corbihelia

Hacrobia

Haptophytina

Retaria

Rhizaria

Cercozoa

Harosa

Heterokonta

Halvaria
CHROMISTA

Alveolata

CORTICATA

Glaucophyta

PLANTAE

Streptophyta

Viridiplantae

Chlorophyta

Centrohelea CHROMISTA

Loukozoa

Metamonada

Malawimonadea

Diphymonadea

Planomonadea

Apusomonadida

Parabasalia

Sulcozoa

Breviatea

ANIMALIA

opisthokonts

Choanoflagellatea

Choanozoa

Filosporidia

Cristidiscoidea

FUNGI

Chytridiomycota

Zygomycota

Dikaryota: higher fungi

Allomycota

Tubulinea

Conosa

Amoebozoa

Discosea

Cutosea

neokaryotes
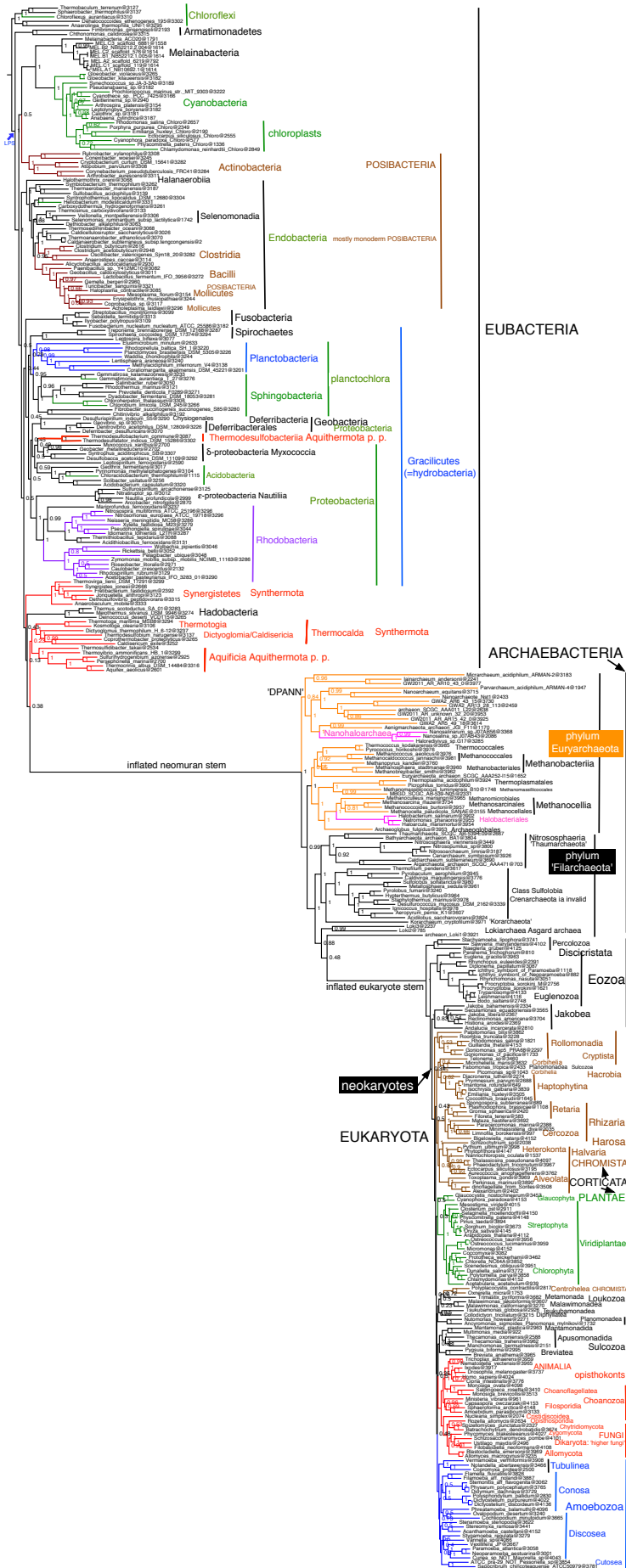
EUKARYOTA

0.9

**Figure S14. Site-homogeneous universal 3-domain RAxML PROTGAMMALGF (4 gamma rates) tree for 26 ribosomal proteins from 143 eukaryotes, 60 archaebacteria, and 151 eubacteria representing all the most divergent lineages.**
Support values are percentages for 100 pseudoreplicates. See next page:

Chloroflexi

Armatimonadetes

Melainabacteria

Cyanobacteria

chloroplasts

Mollicutes
Bacillia

Mollicutes

Clostridiales sensu stricto

Endobacteria
mostly monoderm POSIBACTERIA

Selenomonadia

Thermoanaerobacterales

Heliobacteriales

Sulfobacillales
Halanaerobia

EUBACTERIA

Aquithermota

Thermocalda    Synthermota
Thermotogia
Dictyoglomia/Caldisericia

Synergistetes    Synthermota

Actinobacteria    POSIBACTERIA p. p.

Hadobacteria

Fusobacteria

Rhodobacteria

ε-proteobacteria Nautiliia
Geobacteria

Deferribacterales

Acidobacteria
Nitrospiria

δ-proteobacteria Myxococcia

Spirochaetes

Planctobacteria    planctochlora

Sphingobacteria

Gracilicutes
(=hydrobacteria)

Proteobacteria

inflated neomuran stem

ARCHAEBACTERIA

'Nanohaloarchaea'

'DPANN'

Halobacteriales
Methanocella

Methanosarcinales
Archaeoglobales

Methanomassiliicoccales
Thermoplasmatales

phylum
Euryarchaeota

Methanobacteriales
Thermococcales
Methanococcales

Nitrososphaeria
'Thaumarchaeota'

phylum
'Filarchaeota'

Sulfolobia

Lokiarchaea Asgard archaea
'Korarchaeota'

Percolozoa

inflated eukaryote stem

Eozoa
Euglenozoa

Discicristata

Jakobea

CHROMISTA
Rollomonadia    Cryptista
CORTICATA
Tsukubamonadea
Malawimonadea    Metamonada    Loukozoa
Centrohelea
Planomonadea    Sulcozoa

Corbihelia    Hacrobia

Haptophytina    CHROMISTA

neokaryotes

Streptophyta

Viridiplantae

Chlorophyta

Glaucophyta    PLANTAE

CORTICATA

Alveolata
Halvaria
Heterokonta

CHROMISTA
EUKARYOTA
Retaria    Rhizaria
Cercozoa    Harosa

Chytridiomycota
Zygomycota
Dikaryota: 'higher fungi'
FUNGI
Allomycota
Filosporidia    Choanozoa

ANIMALIA    opisthokonts
Choanoflagellatea    Choanozoa

Breviatea    Sulcozoa
Apusomonadida

Mantamonadida

Tubulinea
Centramoebida

Himatismenida

Conosa
Amoebozoa

Flabellinia

Cutosea

0.5

**Fig. S15. Site-heterogeneous PhyloBayes CAT-GTR tree for 26 ribosomal proteins from 156 eubacteria representing all the most divergent lineages with cultivated repesentatives plus Melainabacteria, chloroplasts and 5 shorter-branch mitochondria.**

Support values are posterior probabilities. Two chains were summed (40,782 trees) after removing 5725 pre-log-likelihood plateau trees as burnin; despite clear plateauing a few deep-branching topological differences (PP 0.5 or less only) remained between the two chains, making maxdiff 1. The corresponding ML tree placed mitochondria one node lower as sister to all α-proteobacteria. None grouped them with the long-branch rickettsias.

**Figure S16.** Site-homogeneous RAxML-PROTGAMMALGF (4 gamma rates) tree of 305 prokaryotic structural maintenace of chromosome (SMC) proteins using 448 amino acid positions. See next page:
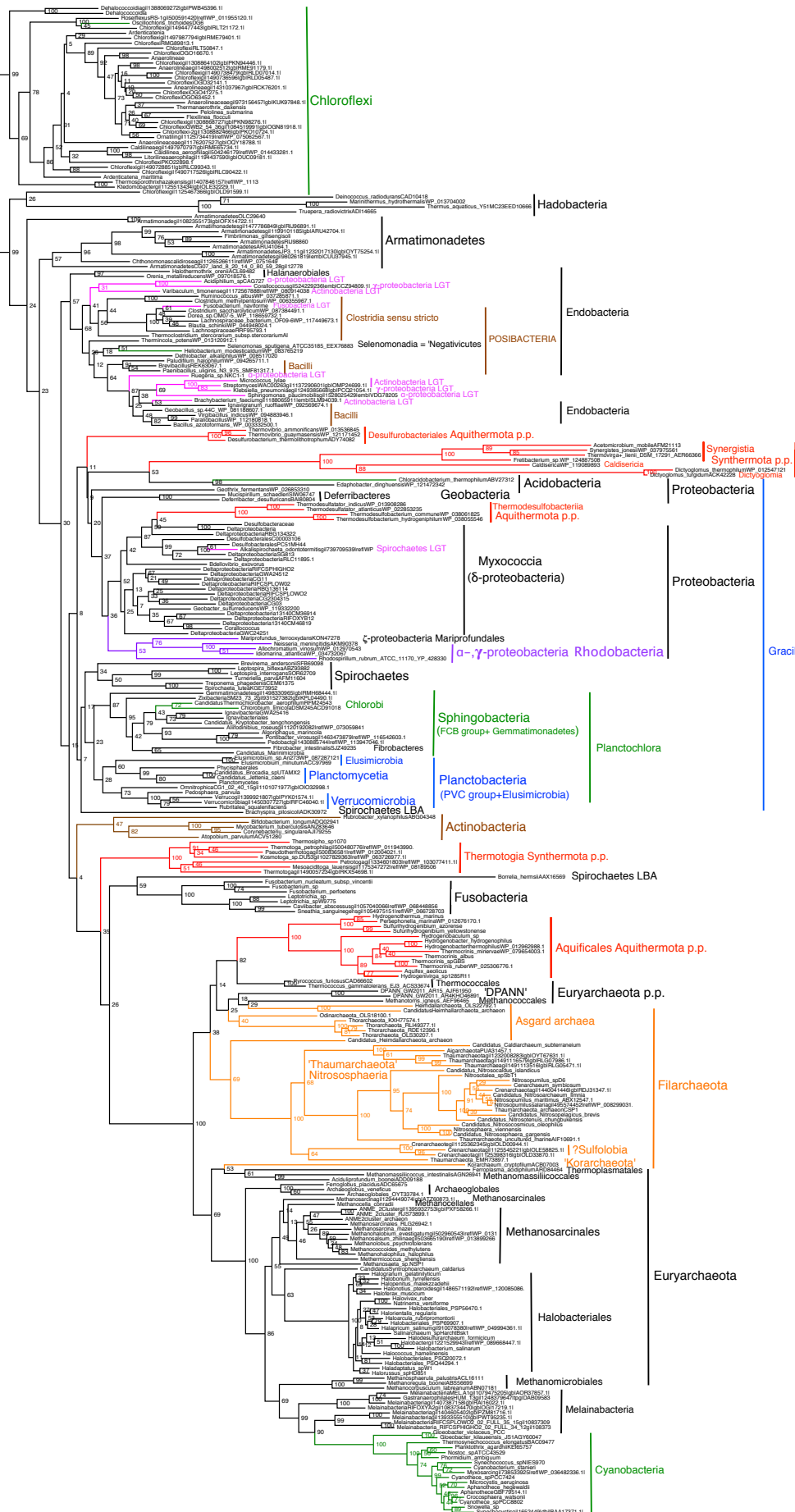
**Figure S16**

**Figure S17. Site-heterogeneous PhyloBayes CAT-GTR tree of 321 largely prokaryotic structural maintenace of chromosome (SMC) proteins using 448 amino acid positions.** maxdiff 0.1246. Four of the five eukaryote sequences group weakly with a relatively short-branch arsgard seuce whereas the human sequence groups with a longer-branch lokiarchaeote sequence. Some other Asgard seunces have even longer branches and group weakly with *Korarchaeum* and strongly with an almiost ceratinly artefactual grouping of Sulfolobia. *Methanopyrus* and the Synthermote eubacterium *Coprothermobacter*, the longest branch on the whole tree.
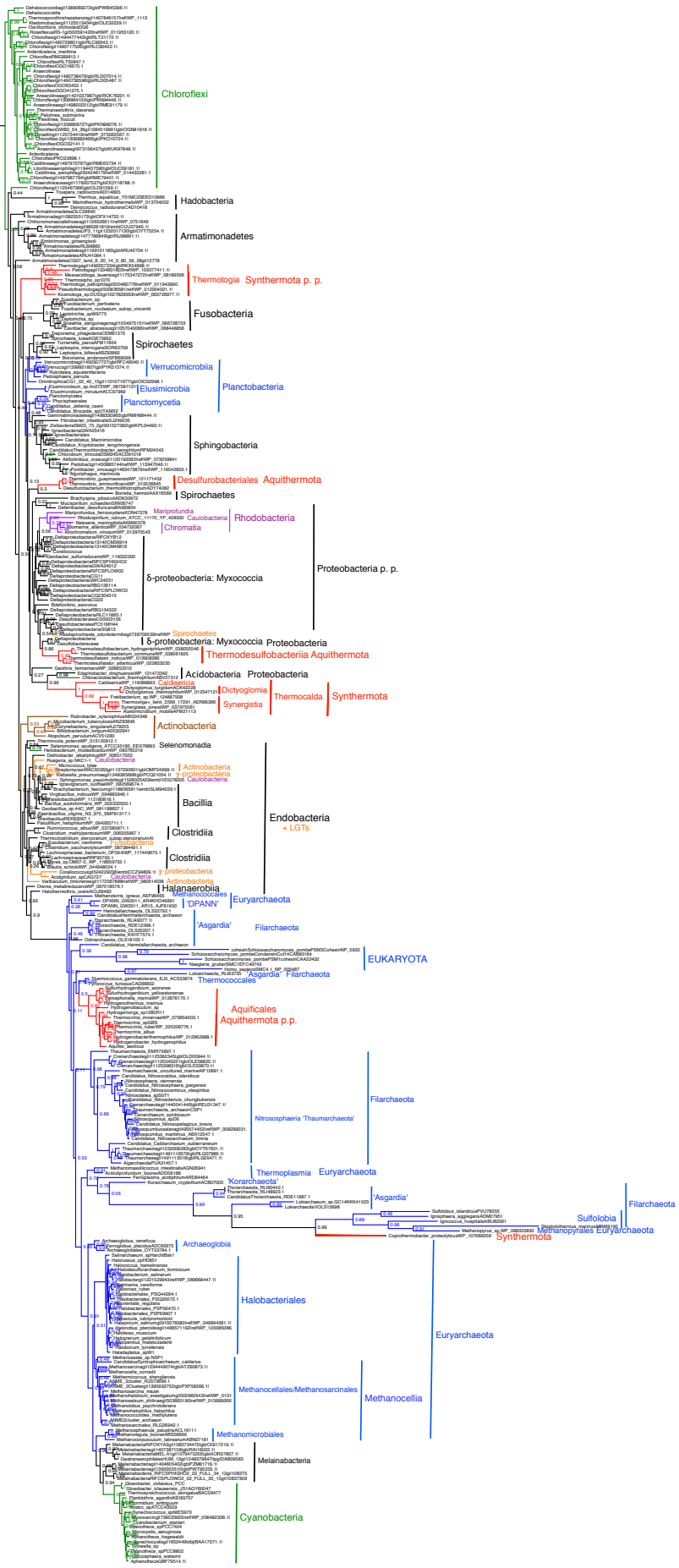
**Figure S17**

**Ribosomal protein (RP) alignments in fasta format as a zip files**
1. These 51 **.fasta files** comprise the regions for 51 RPs after trimming as used for our trees, including all trees shown directly on Figs 3-10 plus those by other methods whose support values are included on them and include over 414 taxa, some excluded from our final analyses.

To enable readers to identify those taxa for which our alignment database uses abbreviated names or interim names, we also include the SCAFOS taxonomy (OTU) file as a key to translate from them to the 354 final names on these trees:
OTU354Arch+Eub+Euk_NoAtrichosa_28July2016.txt

2. **SMC proteins.** 305 complete SMC sequences are in GenBank format (.gb) plus a mask (at the beginning of file) in which 1s mark the 448 included positions and 0s mark those excluded from our analyses.

**Tree files in Newick format in a zip file with separate subfolders for each of the nine figures.**
These comprise all 22 trees used for Figs 3-10 and 12. Their titles indicate the method used (PoiPB means PhyloBayes Poisson; PB means PhyloBayes CAT-GTRGamma)