# Supplemental Material to:

Paweł Mackiewicz, Andrzej Bodył, Krzysztof Moszczyński

## The case of horizontal gene transfer from bacteria to the peculiar dinoflagellate plastid genome

# Materials and methods

## Collection of sequences and preparation of alignments

To gather all potential homologs to Ycf16 and Ycf24 encoded on the *Ceratium horridum* AF490364 minicircle as well as the sequences of Rpl28, Rpl33, and the potential product of the unannotated open reading frame (FtsY) encoded on the *Pyrocystis lunula* AF490367 minicircle, we carried out comprehensive PSI-BLAST or BLAST searches of non-redundant protein and expressed sequence tag databases in NCBI GenBank (http://blast.ncbi.nlm.nih.gov), Dragonblast (http://dbdata.rutgers.edu/dragon), and the expressed sequence tag database of *Alexandrium tamarense*.[1] The identified sequences were verified by local searches of the Conserved Domain Database[2] for the presence of appropriate domains. Each of the resulting sets of homologs consisted of several thousand sequences. Initial alignments were performed in MAFFT using the slow but accurate algorithm L-INS-i with 1,000 cycles of iterative refinement[3] and edited manually in JalView.[4] Incomplete or fragmentary sequences were excluded from further analyses. Ultimately, more than 100 sequences in each of five protein sets were selected using T-Coffee[5] to remove redundancy from the datasets and to include representatives from various prokaryotic and eukaryotic groups. Final alignments were obtained in T-Coffee using profile information (PSI-Coffee) and combining the output of many alignment methods (M-Coffee).

## Phylogenetic analyses

Phylogenetic trees were inferred by six approaches using five programs: PhyloBayes 3.3e (ref. 6), MrBayes 3.2.1 (ref. 7), TreeFinder[8], PhyML-Structure[9], and morePhyML 1.14 (ref. 10) based on PhyML 3.0 (ref. 11). In the PhyloBayes analyses, we applied two substitution models for all alignment sets, LG+$\Gamma$(5) and CAT Poisson+$\Gamma$(5), with the number of components, weights, and profiles inferred from the data. Two independent Markov chains were run through 100,000 cycles (for the Ycf24 set) or 200,000 cycles (for the other sets) with the former model and through 1,000,000 cycles with the latter model. A posterior consensus was calculated from the last 10,000-500,000 trees from each chain after obtaining convergence and good or acceptable runs. In the MrBayes approach, we assumed the mixed+I+$\Gamma$(5) model for all alignment sets to sample appropriate models across the substitution model space in the Bayesian MCMC analysis itself, avoiding the need for a priori model testing. In this analysis, two independent runs starting from random trees were applied, each using eight Markov chains with 40,000,000 generations (for the Rpl28 and Rpl33 sets)

or four Markov chains with 20,000,000 generations (for the other sets). Trees were sampled every 100 generations; to calculate a posterior consensus, we selected trees from the last 5,868,000-26,644,000 generations that reached stationary phase and convergence (the standard deviation of split frequencies stabilized and was less than 0.01).

In the TreeFinder approach, we applied appropriate substitution models that were chosen according to the Propose Model module in this program assuming optimized frequencies of amino acids, whereas the models used in (more)PhyML were selected according to ProtTest 3.2 (ref. 12) assuming optimization of models, branches, and topology of the tree (Table 1S). Search depth was set to 2 in TreeFinder, and the best heuristic search algorithms, NNI and SPR, in (more)PhyML were applied. Edge support was assessed by the bootstrap analysis with 1,000 replicates in each of these two programs. Additionally, we applied the Local Rearrangements-Expected Likelihood Weights method in TreeFinder and the approximate likelihood ratio test (aLRT) based on a Shimodaira-Hasegawa-like procedure in morePhyML.[13] In PhyML-Structure, we used the EX_EHO+$\Gamma$(5) substitution model for all alignment sets, whereas edge support was calculated by aLRT based on the $\chi^2$ test and a Shimodaira-Hasegawa-like procedure. The minimum of these two aLRT support values is shown at selected nodes in the presented trees (Fig. 1S).

Topology tests with 10,000,000 replicates were performed in Consel v0.20 (ref. 14) to compare trees obtained in PhyloBayes under LG+$\Gamma$(5) with alternative topologies that assumed different positions of minicircle *C. horridum* and *P. lunula* sequences (Fig. 1S). Site-wise log-likelihoods for the analyzed trees were calculated in PhyML under the best fitted substitution models found in ProtTest.

Table 1S. Applied substitution models in the analyzed alignment sets.

| Alignment set | PhyloBayes | MrBayes | TreeFinder | (more)PhyML | PhyML-Structure |
|---|---|---|---|---|---|
| FtsY | LG+$\Gamma$(5), CAT Poisson+$\Gamma$(5) | mixed+I+$\Gamma$(5) | LG+F+I+$\Gamma$(5) | LG+$\Gamma$(5) | EX_EHO+$\Gamma$(5) |
| Rpl28 | LG+$\Gamma$(5), CAT Poisson+$\Gamma$(5) | mixed+I+$\Gamma$(5) | witHIV+I+$\Gamma$(5) | LG+F+$\Gamma$(5) | EX_EHO+$\Gamma$(5) |
| Rpl33 | LG+$\Gamma$(5), CAT Poisson+$\Gamma$(5) | mixed+I+$\Gamma$(5) | MIX+F+$\Gamma$(5) | LG+$\Gamma$(5) | EX_EHO+$\Gamma$(5) |
| Ycf16 | LG+$\Gamma$(5), CAT Poisson+$\Gamma$(5) | mixed+I+$\Gamma$(5) | LG+F+I+$\Gamma$(5) | LG+F+$\Gamma$(5) | EX_EHO+$\Gamma$(5) |
| Ycf24 | LG+$\Gamma$(5), CAT Poisson+$\Gamma$(5) | mixed+I+$\Gamma$(5) | LG+F+$\Gamma$(5) | LG+F+$\Gamma$(5) | EX_EHO+$\Gamma$(5) |

**References**

1. Chan CX, Soares MB, Bonaldo MF, Wisecaver JH, Hackett JD, Anderson DM, et al. Analysis of *Alexandrium tamarense* (Dinophyceae) genes reveals the complex evolutionary history of a microbial eukaryote. J Phycol 2012; 48:1130-42.

2. Marchler-Bauer A, Anderson JB, Cherukuri PF, DeWeese-Scott C, Geer LY, Gwadz M, et al. CDD: a conserved domain database for protein classification. Nucleic Acids Res 2005; 33:D192-6.

3. Katoh K, Toh H. Recent developments in the MAFFT multiple sequence alignment program. Brief Bioinform 2008; 9:286-98.

4. Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ. Jalview version 2: A Multiple Sequence Alignment and Analysis Workbench. Bioinformatics 2009; 25:1189-91.

5. Notredame C, Higgins DG, Heringa J. T-Coffee: A novel method for fast and accurate multiple sequence alignment. J Mol Biol 2000; 302:205-17.

6. Lartillot N, Philippe H. A Bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. Mol Biol Evol 2004; 21:1095-109.

7. Ronquist F, Teslenko M, van der Mark P, Ayres DL, Darling A, et al. MrBayes 3.2: efficient Bayesian phylogenetic inference and model choice across a large model space. Syst Biol 2012; 61:539-42.

8. Jobb G, von Haeseler A, Strimmer K. TREEFINDER: a powerful graphical analysis environment for molecular phylogenetics. BMC Evol Biol 2004; 4:18.

9. Le SQ, Gascuel O. Accounting for solvent accessibility and secondary structure in protein phylogenetics is highly beneficial. Syst Biol 2010; 59:277-87.

10. Criscuolo A. morePhyML: improving the phylogenetic tree space exploration with PhyML 3. Mol Phylogenet Evol 2011; 61:944-8.

11. Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. Syst Biol 2010; 59:307-21.

12. Darriba D, Taboada GL, Doallo R, Posada D. ProtTest 3: fast selection of best-fit models of protein evolution. Bioinformatics 2011; 27:1164-5.

13. Anisimova M, Gascuel O. Approximate likelihood-ratio test for branches: a fast, accurate, and powerful alternative. Syst Biol 2006; 55:539-52.

14. Shimodaira H, Hasegawa M. CONSEL: for assessing the confidence of phylogenetic tree selection. Bioinformatics 2001; 17:1246-7.

**Figure 1S.** Phylogenetic trees for FtsY, Rpl28, Rpl33, Ycf16, and Ycf24 sequences inferred in PhyloBayes under the LG+Γ(5) model. Sequences localized to *Pyrocystis lunula* AF490367 and *Ceratium horridum* AF490364 minicircles appear in bold font. Numbers at nodes (in order) correspond to: posterior probabilities estimated in Phylobayes for the LG+Γ(5) and CAT Poisson+Γ(5) models as well as in MrBayes, the minimum of support values calculated by aLRT based on the χ2 test and a Shimodaira-Hasegawa-like procedure in PhyML-Structure, support values obtained by a Shimodaira-Hasegawa-like procedure in morePhyML, PhyML bootstrap values, Local Rearrangements-Expected Likelihood Weights-support values calculated in TreeFinder, and TreeFinder bootstrap values. Values of the posterior probabilities and bootstrap percentages lower than or equal to 0.50 and 50%, respectively, were omitted or indicated by a dash "-". Tables show the results of topology tests comparing the best topology with alternatives that assume different positions of the minicircle-encoded proteins. Topology test results are: the p-value for the approximately unbiased test (au) calculated from the multiscale bootstrap, the non-parametric bootstrap probability calculated from the multiscale bootstrap (np), the bootstrap probability calculated in the non-multiscale manner (bp), the Bayesian posterior probability calculated by the BIC approximation (pp), and the p-values of the Kishino-Hasegawa test (kh), the Shimodaira-Hasegawa test (sh), the weighted Kishino-Hasegawa test (skh), and the weighted Shimodaira-Hasegawa test (wsh).

# FtsY

**Bacteroidetes**

**Cyanobacteria**

**plastid-containing eukaryotes**

**other Bacteria**

| Topology | au | np | bp | pp | kh | sh | wkh | wsh |
|---|---|---|---|---|---|---|---|---|
| 1 | 7·10⁻⁸ | 5·10⁻⁸ | 0 | 7·10⁻²³ | 9·10⁻⁷ | 9·10⁻⁷ | 9·10⁻⁷ | 2·10⁻⁶ |
| 2 | 2·10⁻⁷ | 5·10⁻⁹ | 0 | 2·10⁻²⁸ | 0 | 0 | 0 | 0 |
| 3 | 2·10⁻⁸ | 2·10⁻⁸ | 0 | 2·10⁻²⁹ | 0 | 0 | 0 | 0 |

Gramella forsetii KT0803 (Flavobacteria, Bacteroidetes)
Zunongwangia profunda SM-A87 (Flavobacteria, Bacteroidetes)
Bergeyella zoohelcum CCUG 30536 (Flavobacteria, Bacteroidetes)
Gillisia limnaea DSM 15749 (Flavobacteria, Bacteroidetes)
Kordia algicida OT-1 (Flavobacteria, Bacteroidetes)
Flavobacterium frigoris PS1 (Flavobacteria, Bacteroidetes)
Bizionia argentinensis JUB59 (Flavobacteria, Bacteroidetes)
Mesoflavibacter zeaxanthinifaciens S86 (Flavobacteria, Bacteroidetes)
Joostella marina DSM 19592 (Flavobacteria, Bacteroidetes)
Formosa sp. AK20 (Flavobacteria, Bacteroidetes)
Leeuwenhoekiella blandensis MED217 (Flavobacteria, Bacteroidetes)
Aquimarina agarilytica ZC1 (Flavobacteria, Bacteroidetes)
Galbibacter sp. ck-I2-15 (Flavobacteria, Bacteroidetes)
Flexibacter litoralis DSM 6794 (Cytophagia, Bacteroidetes)
Croceibacter atlanticus HTCC2559 (Flavobacteria, Bacteroidetes)
Capnocytophaga gingivalis ATCC 33624 (Flavobacteria, Bacteroidetes)
Flavobacteria bacterium BBFL7 (Flavobacteria, Bacteroidetes)
Nonlabens dokdonensis DSW-6 (Flavobacteria, Bacteroidetes)
Polaribacter irgensii 23-P (Flavobacteria, Bacteroidetes)
Aequorivita sublithincola DSM 14238 (Flavobacteria, Bacteroidetes)
Flavobacterium johnsoniae UW101 (Flavobacteria, Bacteroidetes)
Dokdonia donghaensis MED134 (Flavobacteria, Bacteroidetes)
Krokinobacter sp. 4H-3-7-5 (Flavobacteria, Bacteroidetes)
Polaribacter sp. MED152 (Flavobacteria, Bacteroidetes)
Lacinutrix sp. 5H-3-7-4 (Flavobacteria, Bacteroidetes)
Blattabacterium sp. Blattella germanica Bge (Flavobacteria, Bacteroidetes)
Blattabacterium sp. Panesthia angustipennis spadica BPAA (Flavobacteria, Bacteroidetes)
Subdoligranulum variabile DSM 15176 (Clostridia, Firmicutes)
**Pyrocystis lunula (Dinoflagellata, Alveolata)**
Pontibacter sp. BAB1700 (Cytophagia, Bacteroidetes)
Marinilabilia salmonicolor JCM 21150 (Bacteroidia, Bacteroidetes)
Porphyromonas catoniae F0037 (Bacteroidia, Bacteroidetes)
Barnesiella intestinihominis YIT 11860 (Bacteroidia, Bacteroidetes)
Prevotella saccharolytica F0055 (Bacteroidia, Bacteroidetes)
Indibacter alkaliphilus LW1 (Cytophagia, Bacteroidetes)
Bacteroidetes oral taxon 274 F0058 (Bacteroidetes)
Prevotella stercorea DSM 18206 (Bacteroidia, Bacteroidetes)
Candidatus Azobacteroides pseudotrichonymphae CFP2 (Bacteroidia, Bacteroidetes)
Bdellovibrio exovorus JSS (Deltaproteobacteria, Proteobacteria)
Niastella koreensis GR20-10 (Sphingobacteria, Bacteroidetes)
Niabella soli DSM 19437 (Sphingobacteria, Bacteroidetes)
Chitinophaga pinensis DSM 2588 (Sphingobacteria, Bacteroidetes)
Dyadobacter fermentans DSM 18053 (Cytophagia, Bacteroidetes)
Salinibacter ruber DSM 13855 (Bacteroidetes)
Candidatus Cloacamonas acidaminovorans Evry (candidate division WWE1)
Candidatus Amoebophilus asiaticus 5a2 (Bacteroidetes)
Cardinium endosymbiont cEper1 of Encarsia pergandiella (Bacteroidetes)
Haliscomenobacter hydrossis DSM 1100 (Sphingobacteria, Bacteroidetes)
Alistipes shahii WAL 8301 (Bacteroidia, Bacteroidetes)
Saprospira grandis DSM 2844 (Sphingobacteria, Bacteroidetes)
Flavobacteria bacterium MS024-2A (Flavobacteria, Bacteroidetes)
Paulinella chromatophora (Cercozoa, Rhizaria)
Synechococcus sp. RCC307 (Oscillatoriophycideae, Cyanobacteria)
Oscillatoria acuminata PCC 6304 (Oscillatoriophycideae, Cyanobacteria)
Gloeobacter violaceus PCC 7421 (Gloeobacteria, Cyanobacteria)
Pseudanabaena sp. PCC 7367 (Oscillatoriophycideae, Cyanobacteria)
uncultured Flavobacteria bacterium (Flavobacteria, Bacteroidetes)
Acidaminococcus fermentans DSM 20731 (Negativicutes, Firmicutes)
Candidatus Pelagibacter ubique HTCC1062 (Alphaproteobacteria, Proteobacteria)
Mesotoga prima MesG1.Ag.4.2 (Thermotogales, Thermotogae)
Oryza sativa Japonica Group (Streptophyta, Viridiplantae)
Zea mays (Streptophyta, Viridiplantae)
Arabidopsis thaliana (Streptophyta, Viridiplantae)
Vitis vinifera (Streptophyta, Viridiplantae)
Physcomitrella patens patens (Streptophyta, Viridiplantae)
Guillardia theta CCMP2712 (Pyrenomonadales, Cryptophyta)
Coccomyxa subellipsoidea C-169 (Chlorophyta, Viridiplantae)
Chlorella variabilis (Chlorophyta, Viridiplantae)
Chlamydomonas reinhardtii (Chlorophyta, Viridiplantae)
Micromonas pusilla CCMP1545 (Chlorophyta, Viridiplantae)
Phaeodactylum tricornutum CCAP 1055 1 (Bacillariophyta, Stramenopiles)
Thalassiosira oceanica (Bacillariophyta, Stramenopiles)
Aureococcus anophagefferens (Pelagophyceae, Stramenopiles)
Cyanidioschyzon merolae strain 10D (Bangiophyceae, Rhodophyta)
Nannochloropsis gaditana CCMP526 (Eustigmatophyceae, Stramenopiles)
Neorickettsia risticii Illinois (Alphaproteobacteria, Proteobacteria)
Propionibacterium propionicum F0230a (Actinobacteridae, Actinobacteria)
Actinomyces sp. oral taxon 178 F0338 (Actinobacteridae, Actinobacteria)
Ilumatobacter coccineum YM16-304 (Acidimicrobidae, Actinobacteria)
uncultured delta proteobacterium HF0070 10I02 (Deltaproteobacteria, Proteobacteria)
Bacteriovorax marinus SJ (Deltaproteobacteria, Proteobacteria)
Helicobacter cinaedi CCUG 18818 (Epsilonproteobacteria, Proteobacteria)
Zymomonas mobilis pomaceae ATCC 29192 (Alphaproteobacteria, Proteobacteria)
Novosphingobium sp. Rr 2-17 (Alphaproteobacteria, Proteobacteria)
Candidatus Puniceispirillum marinum IMCC1322 (Alphaproteobacteria, Proteobacteria)
Orientia tsutsugamushi Ikeda (Alphaproteobacteria, Proteobacteria)
Asticcacaulis excentricus CB 48 (Alphaproteobacteria, Proteobacteria)
Anaplasma phagocytophilum HZ (Alphaproteobacteria, Proteobacteria)
Candidatus Blochmannia floridanus (Gammaproteobacteria, Proteobacteria)
Pirellula staleyi DSM 6068 (Planctomycetia, Planctomycetes)
Schlesneria paludicola DSM 18645 (Planctomycetia, Planctomycetes)
Truepera radiovictrix DSM 17093 (Deinococci, Deinococcus-Thermus)
Lentisphaera araneosa HTCC2155 (Lentisphaeria, Lentisphaerae)
Anaerolinea thermophila UNI-1 (Anaerolineae, Chloroflexi)
Neisseria meningitidis serogroup C (Betaproteobacteria, Proteobacteria)
Acidobacterium capsulatum ATCC 51196 (Acidobacteriales, Acidobacteria)
Methylobacter tundripaludum SV96 (Gammaproteobacteria, Proteobacteria)
Persephonella marina EX-H1 (Aquificales, Aquificae)
Thermanaerovibrio acidaminovorans DSM 6589 (Synergistia, Synergistetes)
Jonquetella anthropi E3 33 E1 (Synergistia, Synergistetes)
Anaerobaculum mobile DSM 13181 (Synergistia, Synergistetes)
Candidatus Riesia pediculicola USDA (Gammaproteobacteria, Proteobacteria)
Thermodesulfatator indicus DSM 15286 (Thermodesulfobacteriales, Thermodesulfobacteria)
alpha proteobacterium HIMB114 (Alphaproteobacteria, Proteobacteria)
Candidatus Liberibacter americanus PW SP (Alphaproteobacteria, Proteobacteria)
Lawsonia intracellularis PHE MN1-00 (Deltaproteobacteria, Proteobacteria)
Slackia piriformis YIT 12062 (Coriobacteridae, Actinobacteria)
Ureaplasma urealyticum serovar 7 ATCC 27819 (Mollicutes, Tenericutes)
Mycoplasma genitalium G37 (Mollicutes, Tenericutes)
Mycoplasma haemofelis Ohio2 (Mollicutes, Tenericutes)
Mycoplasma mycoides (Mollicutes, Tenericutes)
Mycoplasma hyorhinis MCLD (Mollicutes, Tenericutes)
Mycoplasma canis UF31 (Mollicutes, Tenericutes)
Eubacterium cylindroides T2-87 (Erysipelotrichi, Firmicutes)
Turicibacter sanguinis PC909 (Erysipelotrichi, Firmicutes)
Treponema vincentii ATCC 35580 (Spirochaetales, Spirochaetes)
Treponema pallidum pallidum Nichols (Spirochaetales, Spirochaetes)
Buchnera aphidicola BCc (Gammaproteobacteria, Proteobacteria)
Wigglesworthia glossinidia endosymbiont of Glossina morsitans (Gammaproteobacteria, Proteobacteria)

# Rpl28



**Bacteroidetes**

- Candidatus Azobacteroides pseudotrichonymphae CFP2 (Bacteroidia, Bacteroidetes)
- Parabacteroides distasonis ATCC 8503 (Bacteroidia, Bacteroidetes)
- Bacteroides sp. 1 1 14 (Bacteroidia, Bacteroidetes)
- Prevotella nigrescens ATCC 33563 (Bacteroidia, Bacteroidetes)
- Porphyromonas sp. oral taxon 279 F0450 (Bacteroidia, Bacteroidetes)
- Porphyromonas uenonis 60-3 (Bacteroidia, Bacteroidetes)
- Odoribacter splanchnicus DSM 20712 (Bacteroidia, Bacteroidetes)
- Anaerophaga sp. HS1 (Bacteroidia, Bacteroidetes)
- Blattabacterium sp. Panesthia angustipennis spadica BPAA (Flavobacteria, Bacteroidetes)
- Candidatus Sulcia muelleri GWSS (Flavobacteria, Bacteroidetes)
- Candidatus Uzinura diaspidicola ASNER (Flavobacteria, Bacteroidetes)
- Weeksella virosa DSM 16922 (Flavobacteria, Bacteroidetes)
- Ornithobacterium rhinotracheale DSM 15997 (Flavobacteria, Bacteroidetes)
- Lacinutrix sp. 5H-3-7-4 (Flavobacteria, Bacteroidetes)
- Aquimarina agarilytica ZC1 (Flavobacteria, Bacteroidetes)
- Flavobacteriales bacterium ALC-1 (Flavobacteria, Bacteroidetes)
- Zunongwangia profunda SM-A87 (Flavobacteria, Bacteroidetes)
- Capnocytophaga ochracea DSM 7271 (Flavobacteria, Bacteroidetes)
- Flavobacterium sp. CF136 (Flavobacteria, Bacteroidetes)
- Mesoflavibacter zeaxanthinifaciens S86 (Flavobacteria, Bacteroidetes)
- Robiginitalea biformata HTCC2501 (Flavobacteria, Bacteroidetes)
- Pedobacter agri PB92 (Sphingobacteria, Bacteroidetes)
- Sphingobacterium spiritivorum ATCC 33300 (Sphingobacteria, Bacteroidetes)
- Alistipes putredinis DSM 17216 (Bacteroidia, Bacteroidetes)
- Fluviicola taffensis DSM 16823 (Flavobacteria, Bacteroidetes)
- Salinibacter ruber DSM 13855 (BacteroidetesOrderIncertaesedis, Bacteroidetes)
- Rhodothermus marinus DSM 4252 (BacteroidetesOrderIncertaesedis, Bacteroidetes)
- Chitinophaga pinensis DSM 2588 (Sphingobacteria, Bacteroidetes)
- Niastella koreensis GR20-10 (Sphingobacteria, Bacteroidetes)
- Cesiribacter andamanensis AMV16 (Cytophagia, Bacteroidetes)
- Fibrisoma limi BUZ 3 (Cytophagia, Bacteroidetes)
- Runella slithyformis DSM 19594 (Cytophagia, Bacteroidetes)
- Microscilla marina ATCC 23134 (Cytophagia, Bacteroidetes)
- Emticicia oligotrophica DSM 17448 (Cytophagia, Bacteroidetes)
- Cytophaga hutchinsonii ATCC 33406 (Cytophagia, Bacteroidetes)
- Belliella baltica DSM 15883 (Cytophagia, Bacteroidetes)
- Nitritalea halalkaliphila LW7 (Cytophagia, Bacteroidetes)
- Echinicola vietnamensis DSM 17526 (Cytophagia, Bacteroidetes)
- Indibacter alkaliphilus LW1 (Cytophagia, Bacteroidetes)
- Mariniradius saccharolyticus AK6 (Cytophagia, Bacteroidetes)
- Cyclobacterium marinum DSM 745 (Cytophagia, Bacteroidetes)
- Candidatus Amoebophilus asiaticus 5a2 ( Bacteroidetes)
- Pontibacter sp. BAB1700 (Cytophagia, Bacteroidetes)
- Fulvivirga imtechensis AK7 (Cytophagia, Bacteroidetes)
- Marivirga tractuosa DSM 4126 (Cytophagia, Bacteroidetes)
- **Pyrocystis lunula (Dinoflagellata, Alveolata)**
- Haliscomenobacter hydrossis DSM 1100 (Sphingobacteria, Bacteroidetes)
- Saprospira grandis Lewin (Sphingobacteria, Bacteroidetes)

**Cyanobacteria and plastidic version of Rpl28**

- Emiliania huxleyi (Isochrysidales, Haptophyta)
- Isochrysis galbana (Isochrysidales, Haptophyta)
- Prymnesium parvum (Prymnesiales, Haptophyta)
- Pavlova lutheri (Pavlovales, Haptophyta)
- Aureococcus anophagefferens (Pelagophyceae, Stramenopiles)
- Vitis vinifera (Streptophyta, Viridiplantae)
- Arabidopsis thaliana (Streptophyta, Viridiplantae)
- Selaginella moellendorffii (Streptophyta, Viridiplantae)
- Chlamydomonas reinhardtii (Chlorophyta, Viridiplantae)
- Bigelowiella natans (Cercozoa, Rhizaria)
- Micromonas sp. RCC299 (Chlorophyta, Viridiplantae)
- Chlorella variabilis (Chlorophyta, Viridiplantae)
- Paulinella chromatophora (Cercozoa, Rhizaria)
- Prochlorococcus marinus MIT 9312 (Prochlorales, Cyanobacteria)
- Chondrus crispus (Florideophyceae, Rhodophyta)
- Gracilaria tenuistipitata liui (Florideophyceae, Rhodophyta)
- Cyanophora paradoxa (Cyanophoraceae, Glaucophyta)
- Galdieria sulphuraria (Bangiophyceae, Rhodophyta)
- Porphyra umbilicalis (Bangiophyceae, Rhodophyta)
- Cyanidioschyzon merolae strain 10D (Bangiophyceae, Rhodophyta)
- Gloeobacter violaceus PCC 7421 (Gloeobacter, Cyanobacteria)
- Trichodesmium erythraeum IMS101 (Oscillatoriophycideae, Cyanobacteria)
- Pseudanabaena sp. PCC 7367 (Oscillatoriophycideae, Cyanobacteria)
- Synechococcus sp. JA-3-3Ab (Oscillatoriophycideae, Cyanobacteria)
- Neospora caninum Liverpool (Apicomplexa, Alveolata)
- Plasmodium vivax Sal-1 (Apicomplexa, Alveolata)
- Alexandrium tamarense (Dinoflagellata, Alveolata)

**other Bacteria**

- Buchnera aphidicola BCc (Gammaproteobacteria, Proteobacteria)
- Candidatus Blochmannia vafer BVAF (Gammaproteobacteria, Proteobacteria)
- Spirochaeta africana DSM 8902 (Spirochaetales, Spirochaetes)
- Chlorobium phaeobacteroides DSM 266 (Chlorobia, Chlorobi)
- Chlamydophila abortus S26/3 (Chlamydiales, Chlamydiae)
- Methylacidiphilum infernorum V4 (Verrucomicrobia)
- Verrucomicrobium spinosum DSM 4136 (Verrucomicrobiae, Verrucomicrobia)
- Nocardioidaceae bacterium Broad-1 (Actinobacteridae, Actinobacteria)
- Tropheryma whipplei Twist (Actinobacteridae, Actinobacteria)
- Mycobacterium tuberculosis CDC1551 (Actinobacteridae, Actinobacteria)

**mitochondrial version of Rpl28**

- Arabidopsis thaliana (Streptophyta, Viridiplantae)
- Oryza sativa Japonica Group (Streptophyta, Viridiplantae)
- Vitis vinifera (Streptophyta, Viridiplantae)
- Physcomitrella patens patens (Streptophyta, Viridiplantae)
- Selaginella moellendorffii (Streptophyta, Viridiplantae)
- Bathycoccus prasinos (Chlorophyta, Viridiplantae)
- Ostreococcus lucimarinus CCE9901 (Chlorophyta, Viridiplantae)
- Chlorella variabilis (Chlorophyta, Viridiplantae)
- Coccomyxa subellipsoidea C-169 (Chlorophyta, Viridiplantae)
- Phytophthora sojae (Oomycetes, Stramenopiles)
- Blastocystis hominis (Blastocystis, Stramenopiles)
- Phaeodactylum tricornutum CCAP 1055/1 (Bacillariophyta, Stramenopiles)
- Thalassiosira pseudonana CCMP1335 (Bacillariophyta, Stramenopiles)
- Ectocarpus siliculosus (PXclade, Stramenopiles)
- Babesia equi (Apicomplexa, Alveolata)
- Theileria parva strain Muguga (Apicomplexa, Alveolata)
- Plasmodium knowlesi strain H (Apicomplexa, Alveolata)
- Neospora caninum Liverpool (Apicomplexa, Alveolata)
- Toxoplasma gondii ME49 (Apicomplexa, Alveolata)
- Alexandrium tamarense (Dinoflagellata, Alveolata)
- Perkinsus marinus ATCC 50983 (Perkinsea, Alveolata)
- Paramecium tetraurelia strain d4-2 (Ciliophora, Alveolata)
- Tetrahymena thermophila (Ciliophora, Alveolata)
- Naegleria gruberi strain NEG-M (Schizopyrenida, Heterolobosea)
- Cyanidioschyzon merolae strain 10D (Bangiophyceae, Rhodophyta)
- Capsaspora owczarzaki ATCC 30864 (Capsaspora, Ichthyosporea)
- Guillardia theta CCMP2712 (Pyrenomonadales, Cryptophyta)
- Dothistroma septosporum NZE10 (Dikarya, Fungi)
- Podospora anserina S mat+ (Dikarya, Fungi)
- Lachancea thermotolerans CBS 6340 (Dikarya, Fungi)
- Pneumocystis jirovecii (Dikarya, Fungi)
- Pseudozyma antarctica T-34 (Dikarya, Fungi)
- Ceriporiopsis subvermispora B (Dikarya, Fungi)
- Mixia osmundae IAM 14324 (Dikarya, Fungi)
- Rhizoctonia solani AG-1 IA (Dikarya, Fungi)
- Coprinopsis cinerea okayama7#130 (Dikarya, Fungi)

**α-Proteobacteria**

- Micavibrio aeruginosavorus EPB (Alphaproteobacteria, Proteobacteria)
- Anaplasma marginale Puerto Rico (Alphaproteobacteria, Proteobacteria)
- Rickettsia typhi Wilmington (Alphaproteobacteria, Proteobacteria)
- Candidatus Hodgkinia cicadicola Dsem (Alphaproteobacteria, Proteobacteria)

**other Bacteria**

- Aquifex aeolicus VF5 (Aquificales, Aquificae)
- Coraliomargarita akajimensis DSM 45221 (Opitutae, Verrucomicrobia)
- Isosphaera pallida ATCC 43644 (Planctomycetia, Planctomycetes)
- Marinithermus hydrothermalis DSM 14884 (Deinococci, Deinococcus-Thermus)
- Desulfovibrio alaskensis G20 (Deltaproteobacteria, Proteobacteria)
- Coprothermobacter proteolyticus DSM 5265 (Clostridia, Firmicutes)
- Atopobium vaginae PB189-T1-4 (Coriobacteridae, Actinobacteria)
- Propionibacterium acnes SK182B-JCVI (Actinobacteridae, Actinobacteria)
- Desulfomonile tiedjei DSM 6799 (Deltaproteobacteria, Proteobacteria)
- Treponema denticola ASLM (Spirochaetales, Spirochaetes)
- Treponema saccharophilum DSM 2985 (Spirochaetales, Spirochaetes)
- Elusimicrobium minutum Pei191 (Elusimicrobia, Elusimicrobia)
- Helcococcus kunzii ATCC 51366 (Clostridia, Firmicutes)
- Gardnerella vaginalis 1500E (Actinobacteridae, Actinobacteria)
- Candidatus Solibacter usitatus Ellin6076 (Solibacteres, Acidobacteria)
- Paenibacillus terrae HPL-003 (Bacilli, Firmicutes)
- Candidatus Arthromitus sp. SFB-mouse-Yit (Clostridia, Firmicutes)
- Clostridium botulinum CFSAN001628 (Clostridia, Firmicutes)
- Thermoanaerobacterium xylanolyticum LX-11 (Clostridia, Firmicutes)
- Fusobacterium gonidiaformans ATCC 25563 (Fusobacteriales, Fusobacteria)
- Sebaldella termitidis ATCC 33386 (Fusobacteriales, Fusobacteria)
- Oscillochloris trichoides DG-6 (Chloroflexales, Chloroflexi)
- Ktedonobacter racemifer DSM 44963 (Ktedonobacteria, Chloroflexi)
- Caminibacter mediatlanticus TB-2 (Epsilonproteobacteria, Proteobacteria)
- Mycoplasma hominis ATCC 23114 (Mollicutes, Tenericutes)
- Mycoplasma penetrans HF-2 (Mollicutes, Tenericutes)
- Erysipelotrichaceae bacterium 5 2 54FAA (Erysipelotrichi, Firmicutes)
- Oscillibacter valericigenes Sjm18-20 (Clostridia, Firmicutes)
- Helicobacter pylori Hp A-9 (Epsilonproteobacteria, Proteobacteria)

| Topology | au | np | bp | pp | kh | sh | wkh | wsh |
|---|---|---|---|---|---|---|---|---|
| 1 | $1 \cdot 10^{-5}$ | $1 \cdot 10^{-5}$ | $1 \cdot 10^{-5}$ | $4 \cdot 10^{-23}$ | 0.001 | 0.009 | 0.001 | 0.001 |
| 2 | $2 \cdot 10^{-69}$ | $7 \cdot 10^{-25}$ | 0 | $1 \cdot 10^{-43}$ | 0 | 0 | 0 | 0 |
| 3 | $5 \cdot 10^{-12}$ | $3 \cdot 10^{-10}$ | 0 | $8 \cdot 10^{-53}$ | 0 | 0 | 0 | 0 |

**RpI33**

| Topology | au | np | bp | pp | kh | sh | wkh | wsh |
|----------|----|----|----|----|----|----|-----|-----|
| 1 | 5·10⁻⁶ | 3·10⁻⁶ | 4·10⁻⁶ | 4·10⁻²⁵ | 2·10⁻⁵ | 4·10⁻⁵ | 2·10⁻⁵ | 6·10⁻⁵ |
| 2 | 6·10⁻⁸ | 6·10⁻⁸ | 0 | 5·10⁻²⁶ | 7·10⁻⁶ | 7·10⁻⁵ | 7·10⁻⁶ | 1·10⁻⁵ |
| 3 | 4·10⁻¹⁰ | 2·10⁻⁹ | 0 | 2·10⁻⁴⁵ | 0 | 0 | 0 | 0 |
| 4 | 5·10⁻⁵⁴ | 2·10⁻²⁰ | 0 | 2·10⁻⁵⁰ | 0 | 0 | 0 | 0 |

# Ycf16



| Topology | au | np | bp | pp | kh | sh | wkh | wsh |
|----------|------|------|-----|------|-----|-----|------|------|
| 1 | $3 \cdot 10^{-108}$ | $3 \cdot 10^{-30}$ | 0 | $1 \cdot 10^{-80}$ | 0 | 0 | 0 | 0 |
| 2 | $1 \cdot 10^{-36}$ | $7 \cdot 10^{-18}$ | 0 | $2 \cdot 10^{-84}$ | 0 | 0 | 0 | 0 |
| 3 | $3 \cdot 10^{-46}$ | $3 \cdot 10^{-20}$ | 0 | $2 \cdot 10^{-87}$ | 0 | 0 | 0 | 0 |

**Bacteroidetes**

**Cyanobacteria**

**plastid-containing eukaryotes**

**other Bacteria and Archaea**

# Ycf24

**Bacteroidetes**

**Cyanobacteria**

**plastid-containing eukaryotes**

**other Bacteria and Archaea**

| Topology | au | np | bp | pp | kh | sh | wkh | wsh |
|---|---|---|---|---|---|---|---|---|
| 1 | $3\cdot10^{-55}$ | $7\cdot10^{-22}$ | 0 | $2\cdot10^{-90}$ | 0 | 0 | 0 | 0 |
| 2 | $7\cdot10^{-102}$ | $3\cdot10^{-29}$ | 0 | $5\cdot10^{-97}$ | 0 | 0 | 0 | 0 |