# 'Geoarchaeon NAG1' is a deeply-rooting lineage of the archaeal order Thermoproteales rather than a new phylum

Lionel Guy, Anja Spang, Jimmy H. Saw and Thijs J. G. Ettema*

# Supplementary Methods

A set of 57 conserved proteins was used as the basis of this study. The set is fully described in Guy *et al.* (2014). Briefly, archaeal clusters of orthologous genes (arCOGs) (Wolf et al 2012) were filtered, and only clusters present in 90% of a set of 58 representative archaea were retained. Homologs to these clusters were searched for in 10 representative bacteria, as well as in 8 novel archaeal genomes. Clusters that were not present in at least 3 bacterial genomes were discarded, leading to a set of 67 clusters. Paralogs and split proteins were curated manually. A discordance filter was then applied, following a method described elsewhere (Guy et al 2014, Viklund et al 2012), in an effort to remove horizontally transferred genes. Briefly, an individual gene tree is build for each cluster, and each tree is compared to all others, measuring the amount of shared, well-supported bipartitions. Trees are then ranked, with the ones having the least common bipartitions with other trees being the most discordant. Fifteen percent of the most phylogenetically discordant proteins were discarded, yielding to a set of 57 clusters. Of these, 32 are ribosomal protein genes. A 33rd cluster of ribosomal proteins (S14) was added to the set, as it was described to be universal previously (Lecompte et al 2002).

For each cluster, proteins were aligned with mafft-linsi v6.847b (Katoh and Toh 2008), and columns with 50% or more gaps were removed. Alignments were visually inspected for misaligned regions. Maximum likelihood phylogenies were inferred from individual alignments and concatenates with RAxML 7.9.5 (Stamatakis 2006) under CATLG model with 100 non-parametric bootstraps. Bayesian phylogenies were run with PhyloBayes

30    MPI 1.4f (Lartillot et al 2013) under CAT-Poisson model with a discrete gamma distribution

31    of rates across sites. Four chains were run in parallel for ~12,000 generations, discarding

32    the first 8000, and sampling every 50 trees. Convergence was assessed by plotting the

33    variation of key statistics (log-likelihood, length of the tree, alpha parameter and number of

34    modes) across generations, and verifying that the variation inside one chain was greater

35    than across chains. In addition, the consensus trees for each chain were compared to

36    ensure that they were similar and that the overall topology as shown in Figure 1 was

37    identical in all four consensus trees. Trees were visualized with FigTree 1.4.0  (Andrew

38    Rambaut, University of Edinburgh).

39          After alignment and concatenation, a $\chi^2$-test (Viklund et al 2012) was applied to find

40    the most and least compositionally sites. Alignments were built by adding increasing

41    amounts (20 to 100%, by 10% increment) of the most and least biased sites. For each of

42    these alignments, 100 non-parametric bootstraps were inferred as mentioned above. In

43    addition, ML phylogenies were inferred for the most and least biased half of the alignment.

44          Small and large ribosomal RNA subunits from 90 archaea and 10 bacteria (used as

45    outgroup) were individually aligned with mafft-linsi v6.847b (Katoh and Toh 2008) and

46    columns consisting of 50% or more gaps were removed using trimAl (Capella-Gutiérrez et al

47    2009). The alignments were then concatenated and maximum likelihood and Bayesian

48    phylogenies were inferred as above. In the Bayesian analysis, four chains were run for

49    ~80,000 generations, discarding the first 40,000, and sampling every 100 tree. Convergence

50    was assessed as above.

51          The occurrence of discussed arCOGs in NAG1 and in SAGs affiliating with

52    Geoarchaeota (AAA471-B05, AAA471-B23, AAA471-C03, AAA471-L13, AAA471-L14,

53    AAA471-O08)(Rinke et al 2013) was determined with PSI-BLAST (Altschul et al 1997) using

54    aligned arCOGs as queries(Wolf et al 2012) (using an inclusion threshold of E-value <

55    0.001). The assignment of each specific protein to a particular arCOG cluster was

56    subsequently checked by BLAST against the Refseq database

57  (http://www.ncbi.nlm.nih.gov/refseq/) to avoid incorrect annotations due to the presence of

58  truncated sequences only.

59
60
## References
62

63  Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W *et al* (1997). Gapped
64  BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic*
65  *Acids Res* **25:** 3389-3402.

66
67  Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T (2009). trimAl: a tool for automated
68  alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25:** 1972-1973.
69
70  Guy L, Saw JH, Ettema TJG (2014). The archaeal legacy of eukaryotes: a phylogenomic
71  perspective. *Cold Spring Harbor Perspectives in Biology* **In press**.
72
73  Katoh K, Toh H (2008). Recent developments in the MAFFT multiple sequence alignment
74  program. *Brief Bioinformatics* **9:** 286-298.
75
76  Kozubal MA, Romine M, Jennings Rd, Jay ZJ, Tringe SG, Rusch DB *et al* (2012).
77  Geoarchaeota: a new candidate phylum in the Archaea from high-temperature acidic iron
78  mats in Yellowstone National Park. *ISME J* **7:** 622–634.
79
80  Lartillot N, Rodrigue N, Stubbs D, Richer J (2013). PhyloBayes MPI. Phylogenetic
81  reconstruction with infinite mixtures of profiles in a parallel environment. *Syst Biol* **62:** 611-
82  615.
83
84  Lecompte O, Ripp R, Thierry J-C, Moras D, Poch O (2002). Comparative analysis of
85  ribosomal proteins in complete genomes: an example of reductive evolution at the domain
86  scale. *Nucleic Acids Res* **30:** 5382-5390.
87
88  Rinke C, Schwientek P, Sczyrba A, Ivanova NN, Anderson IJ, Cheng J-F *et al* (2013).
89  Insights into the Phylogeny and Coding Potential of Microbial Dark Matter. *Nature* **499:** 431-
90  437.
91
92  Spang A, Hatzenpichler R, Brochier-Armanet C, Rattei T, Tischler P, Spieck E *et al* (2010).
93  Distinct gene set in two different lineages of ammonia-oxidizing archaea supports the
94  phylum Thaumarchaeota. *Trends Microbiol* **18:** 331-340.
95
96  Stamatakis A (2006). RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses
97  with thousands of taxa and mixed models. *Bioinformatics* **22:** 2688-2690.
98
99  Viklund J, Ettema TJG, Andersson SGE (2012). Independent genome reduction and
100 phylogenetic reclassification of the oceanic SAR11 clade. *Mol Biol Evol* **29:** 599-615.
101
102 Wolf YI, Makarova KS, Yutin N, Koonin EV (2012). Updated clusters of orthologous genes
103 for Archaea: a complex ancestor of the Archaea and the byways of horizontal gene transfer.
104 *Biol Direct* **7**.
105
106
107