



HAL
open science

Phylogénie et évolution des Archaea, une approche phylogénomique

Celine Petitjean

► **To cite this version:**

Celine Petitjean. Phylogénie et évolution des Archaea, une approche phylogénomique. Sciences agricoles. Université Paris Sud - Paris XI, 2013. Français. NNT : 2013PA112159 . tel-01070633

HAL Id: tel-01070633

<https://theses.hal.science/tel-01070633>

Submitted on 8 Oct 2014

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Thèse de Doctorat de l'Université Paris-Sud

École Doctorale Gènes, Génomes, Cellules (ED 426)

Présentée par :

Céline Petitjean

Pour l'obtention du grade de :

Docteur ès Sciences de l'Université Paris-Sud

Phylogénie et évolution des Archaea, une approche phylogénomique

Thèse soutenue le 27 septembre 2013, devant le jury composé de :

Simonetta GRIBALDO
Vincent DAUBIN
Pierre CAPY
Olivier LESPINET
David MOREIRA
Céline BROCHIER-ARMANET

Rapporteur
Rapporteur
Examineur
Examineur
Co-directeur de thèse
Co-directeur de thèse

Résumé

En 1977, Carl Woese sépare les procaryotes en deux grands groupes en proposant une nouvelle classification basée sur des critères phylogénétiques. Les Archaea deviennent ainsi un domaine à part entière aux côtés des Bacteria et des Eucarya. Depuis, la compréhension de ce nouveau groupe et de ses relations avec les deux autres domaines, essentielles pour comprendre l'évolution ancienne du vivant, est largement passée par l'étude de leur phylogénie. Presque 40 ans de recherche sur les archées ont permis de faire évoluer leur image : de bactéries vivant dans des milieux spécialisés, souvent extrêmes, on est passé à un domaine indépendant, très diversifié aussi bien génétiquement, métaboliquement ou encore écologiquement. Ces dernières années la barre symbolique de cent génomes complets d'archées séquencés a été franchie et, parallèlement, les projets génomiques et métagénomiques sur des groupes peu caractérisés ou de nouvelles lignées de haut rang taxonomique (e.g. Nanohaloarchaea, Thaumarchaeota, ARMAN, Aigarchaeota, groupe MGC, groupe II des Euryarchaeota, etc.) se sont multipliés. Tout ceci apporte un matériel sans précédent pour l'étude de l'histoire évolutive et de la diversité des Archaea. Les protéines ribosomiques ont été utilisées de façon courante pour inférer la position phylogénétique des nouvelles lignées d'Archaea. Néanmoins, les phylogénies résultantes ne sont pas complètement résolues, laissant des interrogations concernant d'importantes relations de parenté. La recherche de nouveaux marqueurs est donc cruciale et c'est dans ce contexte que mon projet de thèse s'inscrit.

À partir de l'analyse des génomes de deux Thaumarchaeota et d'une Aigarchaeota, nous avons identifié 200 protéines conservées et bien représentées dans les différents phyla d'archées. Ces protéines sont impliquées dans de nombreux processus cellulaires, ce qui peut apporter un signal phylogénétique complémentaire à celui des marqueurs de type informationnel utilisés par le passé. En plus de confirmer la plupart des relations phylogénétiques inférées à partir de ces derniers (i.e., protéines ribosomiques et sous unités de l'ARN polymérase), l'analyse phylogénétique de ces nouveaux marqueurs apporte un signal permettant une meilleure résolution de la phylogénie des archées et la clarification de certaines relations jusqu'ici confuses.

Un certain nombre de ces nouveaux marqueurs sont aussi présents chez les bactéries. Les relations entre les grands phyla d'archées restant encore non résolues, nous avons utilisé ces protéines pour essayer de placer la racine de l'arbre des Archaea en utilisant comme groupe extérieur les bactéries. Nous avons ainsi pu identifier 38 protéines, parmi les 200 sélectionnées précédemment, ayant un signal phylogénétique suffisamment fiable pour cette étude, auxquelles nous avons ajouté 32 protéines ribosomiques universelles. L'utilisation conjointe de ces données nous a permis de placer la racine entre les Euryarchaeota, d'une part, et un groupe rassemblant les Thaumarchaeota, les Aigarchaeota, les Korarchaeota et les Crenarchaeota, d'autre part. Ce nouvel éclairage sur l'évolution ancienne des archées nous a amené à proposer une révision de leur taxonomie avec, principalement, la création du nouveau phylum "Proteoarchaeota" contenant les quatre phyla actuels que nous proposons de rétrograder en classes : Thaumarchaea, Aigarchaea, Korarchaea et Crenarchaea.

Finalement, l'analyse des protéines codées dans les trois génomes qui ont servi de point de départ de ma thèse nous a permis de générer une masse considérable de données qui ont révélé des traits particuliers ou encore des histoires évolutives inattendues. Un exemple est l'histoire du complexe formé par la chaperonne DnaK et de ses co-chaperonnes GrpE, DnaJ, et DnaJ-Fer chez les Thaumarchaeota, impliquant plusieurs transferts horizontaux entre les trois domaines du vivant.

Abstract

In 1977, Carl Woese proposed a new classification of organisms based on phylogenetic criteria where he divided prokaryotes into two major groups. Thus, Archaea were defined as a new domain, together with Bacteria and Eucarya. Since then, the study of this group and its relationships with the two other domains, essential to understand the early evolution of Life, has been largely done through the investigation of its phylogeny. Almost 40 years of research on the archaea have led to a significant evolution of the knowledge on this group: from considering them as bacteria living in specialized environments, most often extreme ones, to defining them as an independent domain, highly diversified in genetic, metabolic and ecological terms. During the last years, the symbolic barrier of 100 complete archaeal genome sequences has been reached and, simultaneously, many genome projects from poorly-known groups or new high-rank lineages (e.g., Nanohaloarchaea, Thaumarchaeota, ARMAN, Aigarchaeota, MGC, group II Euryarchaeota, etc.) have been launched. All this provides unprecedented information to study the evolutionary history of Archaea. Ribosomal proteins have been used recurrently to infer the phylogenetic position of new archaeal lineages. Nevertheless, the resulting phylogenies are not fully resolved and several important nodes remain uncertain. The identification of new phylogenetic markers is therefore crucial. This represents the framework of my PhD thesis project.

On the basis of the analysis of the genome sequences of two Thaumarchaeota and one Aigarchaeota, we have identified 200 conserved proteins well represented among the different archaeal phyla. These proteins are involved in a number of cellular functions, thus providing a phylogenetic signal complementary to the one obtained from the informational proteins (i.e., ribosomal proteins and RNA polymerase subunits). The phylogenetic analysis of these new markers has led to a better resolution of the archaeal phylogeny, including several relationships that remained unclear.

Several of the new markers are also present in bacteria. Since the relationships among the different archaeal phyla are not yet resolved, we have used those markers to try to place the root of the archaeal phylogeny using the bacterial sequences as outgroup. We have identified 38 proteins among the 200 detected before containing a phylogenetic signal useful for that purpose, to which we have added 32 universal ribosomal proteins. The use of this complete dataset allowed us locating the root between the Euryarchaeota and a large group joining the Thaumarchaeota, Aigarchaeota, Korarchaeota and Crenarchaeota. This new result on the ancient evolutionary history of Archaea has led us to propose a taxonomic revision for this domain, in particular the erection of a new phylum "Proteoarchaeota", containing the current four phyla that we propose to retrograde into classes (Thaumarchaeales, Aigarchaeales, Korarchaeales and Crenarchaeales).

Finally, the analysis of the proteins encoded by the three reference genomes at the origin of this work has generated a large amount of data, which reveals particular traits in certain organisms or unexpected evolutionary histories. One example concerns the evolution in Thaumarchaeota of the protein complex composed of the DnaK chaperon and its co-chaperons GrpE, DnaJ, and DnaJ-Fer, which involves several horizontal gene transfer events among the three domains of Life.

Remerciements

Merci aux membres de mon jury, Simonetta Gribaldo, Vincent Daubin, Pierre Capy et Olivier Lespinet d'avoir accepté d'évaluer mon travail de thèse.

Merci à mes deux directeurs de thèse, David Moreira et Céline Brochier-Armanet de m'avoir fait confiance et de m'avoir permis de mener à bien ce travail.

Merci à Purificación López-García de m'avoir fait confiance et accepté dans son équipe.

Merci à Simonetta Gribaldo de m'avoir permis et poussée à entreprendre cette thèse.

Merci à l'Agence Nationale de la Recherche d'avoir financé ma thèse au travers du projet EVOLDEEP et de l'Investissement d'avenir ANCESTROME pour l'accès à certains clusters de calculs ; à l'Université Paris-sud pour l'accès au cluster de calcul ebio.

Merci à l'École Doctorale GGC de m'avoir acceptée malgré la complexité de mon dossier.

Au début de mon stage de M2, Simonetta m'a dit un jour « tu verras, bientôt le labo sera ta deuxième maison » ; je ne soupçonnais pas à quel point cela serai vrai.

Alors merci aussi à :

Simonetta pour ton encadrement, ta confiance et ta passion pour la science et la phylogénie, et aussi de t'être battue pour moi et de m'avoir poussée à rejoindre finalement David et Céline pour ma thèse.

Alexis pour toute l'aide que tu m'as apportée pendant mon M2 et après, et pour toutes nos discussions. Elie, d'avoir été un point de repère en tant que thésard et pour ton aide dans des moments critiques de fin de stage.

Céline, de m'avoir acceptée en thèse, de m'avoir fait confiance. Merci aussi de m'avoir fait partager ta passion, de ta rigueur scientifique, et de ta bienveillance, particulièrement dans des moments compliqués.

Laura, ma co-thésarde du LCB, pour tout ce qu'on a pu partager, scientifiquement et personnellement ; d'être là tout simplement.

Rym, pour toutes tes questions, ta passion et notre amitié. Rémi, Sandrine, Boyang de l'équipe GEB, Mélodie, ma stagiaire, et aux membres du LCB, de tous nos échanges.

Toute l'équipe microbio de l'ESE, de votre passion scientifique, d'être aussi soudés et de pouvoir partager autant avec vous.

David M., de m'avoir acceptée en thèse, de m'avoir fait confiance, de partager ton incroyable savoir, de la justesse de ton encadrement, de ta bienveillance, et pour ton humour caustique.

Purificación, de ta droiture, de tes coup de gueules, de toutes nos discussions, d'avoir su créer avec David cette équipe, de ta confiance, et particulièrement, de m'avoir poussé à raciner l'arbre des archées. Philippe, pour toute ton aide, de ta capacité à saisir certaines choses, de ton humour, d'être passé à Linux ; et Hélène, de ta présence autour de notre équipe. Bienvenue à Juliette Rusticule. Ludwig et Paola, de votre amitié et de votre présence. Marie, coloc et co-bureau, de ton amitié, d'être la première thésarde de l'équipe microbio et toujours présente. Charles, mon « grand frère » du labo, d'être aussi adorable, de tes talents d'imitateur, et de ta thèse qui m'a bien inspirée. Estelle, co-bureau et voisine, de ton soutien, de ton aide et de ce qu'on a partagé du côté d'Alesia.

Jonathan, pour répondre à tes remerciements, te supporter a été un immense plaisir et une grande richesse, au labo, à la maison et le reste du temps, mais aussi depuis bientôt 10 ans ; de toutes nos discussions enflammées ou calmes, de ton humour de ton soutien et de tout ce pourquoi je ne trouverai pas les mots, merci. Marianne, de ton aide et de ta gentillesse, de tous les CMD. Je n'ai aucun doute sur le fait que tu finiras ta thèse avec panache. Aurélien, je ne doute pas que ta thèse sera une belle aventure.

Julien, mon co-bureau direct, de toutes nos discussions sur la phylogénie et sur la vie, et pour mendeley. Tous les membres du bureau 208 et de la salle T pour nos passionnantes discussions, à 5h et le reste du temps. Pour ceux que je n'ai pas encore nommés, Boris et Yann ; Vincent, passé rapidement. Aux autres membres de l'ESE, si intéressants à côtoyer ; particulièrement Gwendal, Jacqui, Amandine, Hervé, Lucie, Alodie, Martha... c'est un vrai plaisir que de travailler dans cet environnement.

Aux membres du LBBE, croisés trop peu souvent, mais toujours avec plaisir.

Aux secrétaires, Vanessa, Nathalie et Emmanuelle, qui m'ont largement aidée dans mes nombreux problèmes administratifs.

Tous les profs qui m'ont apporté savoirs, savoir-faire et méthodes. Particulièrement mes profs de bio de collège et lycée, Philippe Kachidian, Pascal Hingamp, Keith Dudley, Michel Termier, Daniel Gautheret, Olivier Lespinet, Dominique De Vienne, Fabrice Confalonieri et Pierre Capy de m'avoir encouragé à poursuivre. De même que Patrick Forterre, Eduardo Rocha et Pierre-Henri Gouyon pour les discussions partagées. Emese et Carl, ça a été un plaisir d'enseigner à vos côtés. Pierre, de tes encouragements, de ta disponibilité, de ta bienveillance et de m'avoir accepté au sein de GGC. Olivier et Pascal, mes « parrains », d'avoir accepté ce rôle. Pascal, de ta compréhension et de ta disponibilité quand j'en ai eu besoin.

Marie-Françoise et Patrick, de m'avoir toujours encouragée et soutenue dans mes choix. Aurore et Bérandère, mes merveilleuses petites sœurs, Davi mon cousin adoré, Isabelle et Martine pour être là, toujours.

Mes amis, tous, depuis longtemps et de tous les côtés, qui m'aident à avancer jour après jour. Parmi eux, Claire, David et Doris si importants où que vous soyez ; mes coloc actuels Gregory et Virginie, qui m'avez offert une belle mise au vert pour ma rédaction ; Sarah et Sylvain, toujours là ; et tout ceux que je ne peux pas nommer ici mais qui n'en restent pas moins essentiels. L'AMZ, l'ESP, Cannes et Dragons et Danielle Pauly, & Cie. de m'avoir permis de décompresser et de tenir le coup !

Aux anciens thésards dont j'ai relu les thèses, et à tous les relecteurs de la mienne, Bérandère, David, Jonathan, Laura, Marie, Marianne et Sébastien, MERCI, avec une mention spéciale à Sébastien et Marianne pour en avoir relu la quasi-totalité.

Bestiolus, d'avoir réalisé cette thèse avec moi, sans bug majeur en trois ans.

Sébastien, pour ton incommensurable soutien ces derniers temps, pour ton aide précieuse pour la toute fin, pour ta patience, et pour tout ce que je ne sais comment exprimer. Je ne m'étends pas plus, je mets mon chapeau bleu, et j'arrive.

En hommage posthume à Carl Woese (1928- 2012)

« Tel est le pouvoir des lettres quand seulement l'ordre en est changé »

Lucrèce, *De rerum natura*

Sommaire

RESUME	3
ABSTRACT	4
REMERCIEMENTS	5
SOMMAIRE.....	9
TABLE DES FIGURES	13
INTRODUCTION	15
A. Histoire de la découverte des Archaea.....	15
1. Les trois domaines du vivant	15
a. Un marqueur moléculaire pour la classification bactérienne : l'ARN ribosomique	15
b. Premières phylogénies moléculaires microbiennes et découverte des « Archaeobacteria ».....	16
c. La nouvelle classification du vivant : « Archaea », « Bacteria », « Eucarya ».....	21
2. A partir des années 1990 : Exploration de la diversité archée.	23
a. Etudes environnementales	23
b. Les années 2000 : Séquençage de nombreux génomes complets.....	28
3. Conclusion	30
B. Phylogénie des Archaea.....	31
1. Diversité des Archaea : la phylogénie de l'ARNr SSU	31
a. Les Euryarchaeota.....	33
Hyperthermophiles	33
Mésophiles et psychrophiles	35
Halophiles.....	37
Acidophiles	38
Méthanogènes et méthanotrophes.....	40
b. Les Nanoarchaeota	42
c. Les Crenarchaeota	43
d. Les Thaumarchaeota.....	46
e. Les Aigarchaeota	50
f. Les Korarchaeota	50
g. Conclusion	52

Phylogénie et évolution des Archaea, une approche phylogénomique

2.	Phylogénies moléculaires inférées sur plusieurs marqueurs	53
a.	La phylogénie des archées	53
	Première phylogénie des archées.....	53
	Evolution de la méthanogenèse et implications pour la phylogénie des archées	55
	La position de Nanoarchaeum equitans : nouveau phylum ou artefact ?.....	58
	Nouveaux génomes et nouveaux marqueurs	60
	Nouveaux phyla et phylogénie de référence actuelle	63
b.	La racine de l'arbre des archées et la relation avec les autres domaines	70
3.	Conclusion	75

OBJECTIFS..... 77

Objectif 1 : La recherche de nouveaux marqueurs pour l'inférence de la phylogénie des archées.	77
Objectif 2 : La recherche de la racine de l'arbre des archées grâce à des homologues bactériens.	77

MATERIELS ET METHODES..... 79

A. Analyses phylogénétiques des protéines codées dans les génomes de *C. symbiosum*, *N. maritimus* et '*Ca.*

<i>Caldiarchaeum subterraneum</i>'.....	80
1. Construction d'une banque de données locale	80
2. Génération des phylogénies préliminaires	80
3. Tri et sélection des phylogénies préliminaires.....	82

B. Inférence de la phylogénie des Archaea : méthodes du chapitre 1.....84

1. Analyse des protéines d'intérêt pour l'étude de la phylogénie des Archaea	84
a. Sélection des protéines d'intérêt	84
b. Construction d'une banque de données locale de génomes complets d'Archaea.....	84
c. Construction des jeux de données.....	84
d. Analyse phylogénétique des jeux de données.....	86
Alignement	86
Phylogénies préliminaires.....	86
Répartition taxonomique des séquences	86
Analyse de l'alignement.....	87
Phylogénies individuelles définitives	88
e. Analyse fonctionnelle des 200 nouveaux marqueurs	88
2. Mise à jour des jeux de données de protéines informationnelles.....	88
3. Inférence de la phylogénie globale des Archaea	89
a. Construction des supermatrices	89
b. Désaturation	89
Désaturation par sélection de sites	90

Désaturation par sélection de gènes	91
c. Inférences phylogénétiques	91
C. Racinement de l'arbre des Archaea : méthodes du chapitre 2	93
1. Génération des jeux de données pour le racinement de l'arbre des Archaea	93
a. Sélection des protéines d'intérêt	93
b. Construction d'une banque de données locale de génomes complets de bactéries.....	93
c. Construction des jeux de données pour l'analyse des marqueurs potentiels.....	93
d. Analyse des jeux de données	94
2. Mise à jour des jeux de données de protéines informationnelles pour le racinement de l'arbre des Archaea ..	94
3. Inférence phylogénétique et racinement de l'arbre des Archaea	95
a. Construction des supermatrices	95
b. Désaturation	95
Désaturation par sélection de sites	96
Désaturation par sélection de gènes	96
c. Inférences phylogénétiques	96
 CHAPITRE 1 : LA PHYLOGENIE DES ARCHEES, AU-DELA DES PROTEINES	
INFORMATIONNELLES.....	97
A. Introduction	97
B. Manuscrit de l'article 1 : «Extending the conserved phylogenetic core of Archaea disentangles the evolution of the third domain of Life. »	99
C. Synthèse et éléments de discussion.....	125
 CHAPITRE 2 : POSITIONNEMENT DE LA RACINE DE L'ARBRE DES ARCHAEA.....	129
A. Introduction	129
B. Manuscrit de l'article 2 : « <i>Phylogenomic Analysis Pinpoints the Root of the Domain Archaea and Supports the Foundation of the New Kingdom Proteoarchaeota</i> »	131
C. Synthèse et éléments de discussion.....	155
 CHAPITRE 3 : REPARTITION TAXONOMIQUE ET HISTOIRE EVOLUTIVE DES PROTEINES	
D'ARCHEES : EXEMPLE DU SYSTEME CHAPERONNE DNAK ET DE LA PROTEINE DNAJ-FER.	
.....	157

A. Introduction	157
B. Manuscrit de l'article 3 : « <i>Horizontal gene transfer of a chloroplast DnaJ-Fer protein to Thaumarchaeota and the evolutionary history of the DnaK chaperone system in Archaea</i> »	159
C. Eléments de discussion.....	174
DISCUSSION.....	177
De l'importance de la phylogénie en biologie.	177
De la sélection des données.	178
De l'augmentation de la quantité de données génomiques.	179
De l'intérêt des analyses automatiques.	180
De la multiplication des phyla archéens.	182
Du risque de découplage entre génome et organisme.	183
De la phylogénie moléculaire comme révolution scientifique ?.....	184
CONCLUSION	185
BIBLIOGRAPHIE.....	187
ANNEXES	207

Table des figures

Figure 1. Table représentant les coefficients d'association (SAB) entre 13 représentants des trois royaumes primaires.....	16
Figure 2. Représentation des grandes lignées du vivant selon Fox et collaborateurs en 1980.	18
Figure 3. Premier dendrogramme représentant les relations entre les « archaeobacteria » selon Fox et collaborateurs en 1980.	19
Figure 4. Arbre du vivant représentant les trois domaines Bacteria, Archaea et Eucarya.	22
Figure 5. Phylogénie des Archées montrant les groupes I et II.....	23
Figure 6. Images d'archées non cultivées en FISH.	26
Figure 7. Phylogénie des archées cultivées et non cultivées.	32
Figure 8. Champ hydrothermal d'El Tatio au Chili.	34
Figure 9. Euryarchées hyperthermophiles.	35
Figure 10. Diversité des archées mésophiles.....	36
Figure 11. Archées halophiles.....	38
Figure 12. Mine acide à Rio Tinto en Espagne.	39
Figure 13. Archées acidophiles.	39
Figure 14. Archées méthanogènes et leurs habitats.	41
Figure 15. Nanoarchaeum equitans et Ignicoccus hospitalis, son hôte.	43
Figure 16. Fumeur noir de la dorsale océanique East Pacific Rise.	44
Figure 17. Crenarchées.	46
Figure 18. Thaumarchées.	47
Figure 19. Phylogénie des Archées centrée sur les Thaumarchaeota non cultivées.....	49
Figure 20. 'Candidatus Korarchaeum cryptofilum'.	52
Figure 21. Phylogénies des archées publiées en 2002.	54
Figure 22. Position phylogénétique de M. kandleri.....	56
Figure 23. Nanoarchaeum equitans : une Euryarchaeota.....	59
Figure 24. Phylogénie des archées en 2006.	61
Figure 25. Les Thaumarchaeota : un nouveau phylum?.....	64
Figure 26. Phylogénie des archées en 2011 (référence actuelle).....	67
Tableau 1. Etudes de la phylogénie des archées : génomes et marqueurs utilisés.	70
Tableau 2. Analyses phylogénomiques sur l'arbre du vivant.	72
Tableau 3 : Nombre de protéines conservées à chaque étape de l'analyse des génomes complets de <i>N. maritimus</i> , <i>C. symbiosum</i> et ' <i>Ca. Caldiarchaeum subterraneum</i> '	82
Tableau 4 : Nombre de protéines conservées à chaque étape de l'analyse des marqueurs potentiels depuis les génomes complets de <i>N. maritimus</i> , <i>C. symbiosum</i> et ' <i>Ca. Caldiarchaeum subterraneum</i> '	87
Tableau 5 : Les différentes supermatrices construites pour les analyses phylogénétiques	89
Tableau 6 : Nombre de jeux de données conservés à chaque étape de l'analyse visant à raciner l'arbre des Archaea avec les homologues bactériens.....	94
Tableau 7 : Différentes concaténations utilisées pour les analyses phylogénétiques de racinement de l'arbre des archées.....	95

Introduction

A. Histoire de la découverte des Archaea

1. Les trois domaines du vivant

a. Un marqueur moléculaire pour la classification bactérienne : l'ARN ribosomique

Au début des années 1970, suivant la proposition de Darwin en 1859 (Darwin 1859), classer les organismes vivants selon leur généalogie, puis celle de Zuckerkandl et Pauling en 1965 (Zuckerkandl and Pauling 1965), utiliser les molécules portant l'information génétique pour l'inférence phylogénétique (à savoir les gènes ou les molécules issues de l'expression des gènes (protéines et ARN)), Woese a utilisé les ARN ribosomiques afin de retracer les relations phylogénétiques entre organismes vivants. Des mesures de distances phylogénétiques avaient déjà été construites sur la base de petits polypeptides tels que les hémoglobines ou le cytochrome c, mais pas sur l'ensemble du vivant (Fitch and Margoliash 1967), et à partir de mesures d'hybridation des gènes codant les ARN ribosomiques mais de façon peu précise (Moore and McCarthy 1967; B. Pace and Campbell 1971). En trouvant un moyen d'analyser la séquence assez longue des ARN ribosomiques, Woese a pu appliquer sa méthode d'analyse sur tout organisme vivant, le ribosome et ses composés protéiques et nucléiques étant homologues dans l'ensemble des organismes connus. Pour cela, il a utilisé en premier lieu les molécules d'ARN 5S de la grande sous-unité du ribosome (ARNr 5S) digérées par deux nucléases (nucléase T1 et une nucléase pancréatique). Les oligonucléotides résultants pouvant être séquencés, il a ensuite comparé les séquences et les fréquences de ces oligomères pour définir une distance phylogénétique, « PD value » (« Phylogenetic Distance value »), entre les différents organismes étudiés (Sogin, Sogin, and Woese 1972). Cette distance phylogénétique est un pourcentage correspondant au nombre de changements de bases nécessaires pour passer du catalogue d'oligonucléotides d'une espèce à celui d'une autre espèce. C'est une forme d'approximation du nombre de substitutions entre les séquences de deux organismes. Les premières analyses ont été conduites sur quelques bactéries et ont permis de mettre en évidence l'intérêt de cette méthode. En effet, la distance phylogénétique ainsi mesurée est alors inférieure à 15% entre deux individus d'une même famille, par exemple entre deux entérobactéries, *Escherichia coli* et *Proteus mirabilis* (Sogin, Sogin, and Woese 1972; Carl R Woese et al. 1976). L'intérêt principal de cette méthode était la possibilité de mesurer la distance évolutive entre espèces bactériennes d'une manière objective. Auparavant, seuls des critères phénotypiques (morphologie, tolérance à l'oxygène, capacité à utiliser le glucose, etc.) étaient employés pour

distinguer les espèces bactériennes, mais les microbiologistes étaient bien conscients du fait que les classifications ainsi obtenues n'étaient pas "naturelles", i. e. elles ne reflétaient pas les relations évolutives entre espèces. Ceci avait amené le célèbre microbiologiste américain Roger Stanier à déclarer avec pessimisme que « *the ultimate scientific goal of biological classification cannot be achieved in the case of bacteria* » (Stanier, R. Y., Doudoroff, M. & Adelberg 1963). Le travail de Woese prouva le contraire.

b. Premières phylogénies moléculaires microbiennes et découverte des « Archaeobacteria »

En 1977, Woese et ses collaborateurs analysent l'ARN ribosomique de la petite sous-unité du ribosome (ARNr SSU), 16S chez les bactéries et 18S chez les eucaryotes, et calculent un coefficient d'association entre deux catalogues d'oligonucléotides A et B : le S_{AB} . Il correspond à deux fois le nombre d'oligonucléotides (hexamères ou plus longs) communs entre les deux catalogues sur la somme des oligonucléotides totaux des deux catalogues. Ainsi plus deux séquences sont proches, plus le nombre d'oligonucléotides en commun est important et plus le S_{AB} sera élevé (George E Fox et al. 1977; G. E. Fox, Pechman, and Woese 1977) (Figure 1).

	1	2	3	4	5	6	7	8	9	10	11	12	13
1. <i>Saccharomyces cerevisiae</i> , 18S	—	0.29	0.33	0.05	0.06	0.08	0.09	0.11	0.08	0.11	0.11	0.08	0.08
2. <i>Lemna minor</i> , 18S	0.29	—	0.36	0.10	0.05	0.06	0.10	0.09	0.11	0.10	0.10	0.13	0.07
3. L cell, 18S	0.33	0.36	—	0.06	0.06	0.07	0.07	0.09	0.06	0.10	0.10	0.09	0.07
4. <i>Escherichia coli</i>	0.05	0.10	0.06	—	0.24	0.25	0.28	0.26	0.21	0.11	0.12	0.07	0.12
5. <i>Chlorobium vibrioforme</i>	0.06	0.05	0.06	0.24	—	0.22	0.22	0.20	0.19	0.06	0.07	0.06	0.09
6. <i>Bacillus firmus</i>	0.08	0.06	0.07	0.25	0.22	—	0.34	0.26	0.20	0.11	0.13	0.06	0.12
7. <i>Corynebacterium diphtheriae</i>	0.09	0.10	0.07	0.28	0.22	0.34	—	0.23	0.21	0.12	0.12	0.09	0.10
8. <i>Aphanocapsa</i> 6714	0.11	0.09	0.09	0.26	0.20	0.26	0.23	—	0.31	0.11	0.11	0.10	0.10
9. Chloroplast (<i>Lemna</i>)	0.08	0.11	0.06	0.21	0.19	0.20	0.21	0.31	—	0.14	0.12	0.10	0.12
10. <i>Methanobacterium thermoautotrophicum</i>	0.11	0.10	0.10	0.11	0.06	0.11	0.12	0.11	0.14	—	0.51	0.25	0.30
11. <i>M. ruminantium</i> strain M-1	0.11	0.10	0.10	0.12	0.07	0.13	0.12	0.11	0.12	0.51	—	0.25	0.24
12. <i>Methanobacterium</i> sp., Cariaco-isolate JR-1	0.08	0.13	0.09	0.07	0.06	0.06	0.09	0.10	0.10	0.25	0.25	—	0.32
13. <i>Methanosarcina barkeri</i>	0.08	0.07	0.07	0.12	0.09	0.12	0.10	0.10	0.12	0.30	0.24	0.32	—

Figure 1. Table représentant les coefficients d'association (SAB) entre 13 représentants des trois royaumes primaires. Adapté de l'article de 1977 de Woese et Fox, « Phylogenetic structure of the prokaryotic domain: the primary kingdoms » (C R Woese and Fox 1977). Les organismes 1, 2 et 3 sont des eucaryotes, 4 à 9 des bactéries et de 10 à 13 des archées. Les SAB entre organismes du même groupe sont tous supérieurs à 0,19 (en rouge pour les eucaryotes, bleu pour les bactéries et vert pour les archées) alors qu'entre organismes de deux groupes différents, la valeur maximale est de 0,14.

Les résultats de S_{AB} calculés pour chaque couple d'espèces donnent une matrice de caractères numériques, utilisable pour la construction d'un dendrogramme représentant les relations phylogénétiques entre ces espèces. Après avoir analysé l'ARNr SSU de quelques bactéries

seulement, Fox et Woese décident de conduire une étude à plus large échelle, sur 40 bactéries représentant la diversité disponible à ce moment là (Balch et al. 1977). Le résultat est surprenant pour les auteurs, car il montre une divergence bien plus forte entre les deux « bactéries » méthanogènes (*Methanobacterium ruminantium* et *Methanobacterium thermoautotrophicum*) et le reste des bactéries analysées, qu'entre ces dernières. Ils proposent alors qu'une divergence très ancienne a donné naissance à deux lignées bien distinctes, l'une contenant les bactéries méthanogènes, l'autre contenant toutes les autres bactéries. La même année, Woese propose une nouvelle classification des procaryotes (C R Woese and Fox 1977) en créant la notion de « royaume primaire » ou « royaume », qui se réfère aux premières unités phylogénétiques dans l'arbre du vivant, donc les premiers niveaux de classification à utiliser. Ce nouveau mode de classification se place en opposition aux classifications utilisées à ce moment là, à savoir la dichotomie entre procaryotes et eucaryotes et la classification de Whittaker en cinq règnes (Animaux, Plantes, Champignons, Protistes et Monères), qui ne sont pas basées sur des critères phylogénétiques (Whittaker 1969). Fox et Woese vont définir trois royaumes, les « Eubacteria », les « Archaeobacteria » et les « Urcaryotes ». Pour la première fois, les microbiologistes disposent d'une classification basée sur un critère objectif et phylogénétique, en accord avec la théorie de l'évolution telle que décrite par Darwin dans son livre *On the Origin of Species* (Darwin 1859), et c'est en cela qu'elle revêt une importance toute particulière.

Dans cette première classification les « archaeobactéries » étaient composées uniquement des méthanogènes, puisqu'aucun autre ARNr SSU d'archée n'avait encore été étudié. Néanmoins, Fox et Woese proposent que d'autres bactéries puissent aussi être des « archaeobactéries » : certaines bactéries halophiles extrêmes dont les parois cellulaires ne contiennent pas de peptidoglycane, comme chez les méthanogènes. Au cours des années suivantes, Woese et ses collaborateurs analysent d'autres ARNr SSU, et cela leur permet de confirmer la place d'autres méthanogènes et des halophiles extrêmes au sein des « archaeobactéries » (Magrum, Luehrsen, and Woese 1978); suivent différents « thermoacidophiles », comme *Sulfolobus acidocaldarius* (C R Woese, Magrum, and Fox 1978; Atoji et al. 1980), et les *Thermoplasma*, classés jusqu'ici parmi les mycoplasmes à cause de leur absence de paroi cellulaire mais que l'analyse de l'ARNr SSU de *Thermoplasma acidophilum* permet de les reclasser parmi les archaeobacteria (C R Woese, Magrum, and Fox 1978; C R Woese, Maniloff, and Zablen 1980). Un article de 1980 « *The phylogeny of Prokaryotes* » de Fox et collaborateurs (Atoji et al. 1980) reprend l'ensemble des résultats accumulés pendant les années 70 et propose une phylogénie schématique du vivant ([Figure 2](#)) et plusieurs dendrogrammes pour les « archaeobactéries » ([Figure 3](#)) et différentes familles bactériennes. On peut déjà voir plusieurs des ordres d'archées dans ce dendrogramme, à savoir les *Halobacteriales*, les

Methanomicrobiales, les *Methanococcales*, les *Thermoplasmatales* (représentés par *Thermoplasma acidophilum*) et les Sulfolobales (représentés par *Sulfolobus acidocaldarius*), seuls représentants des futurs Crenarchaeota.

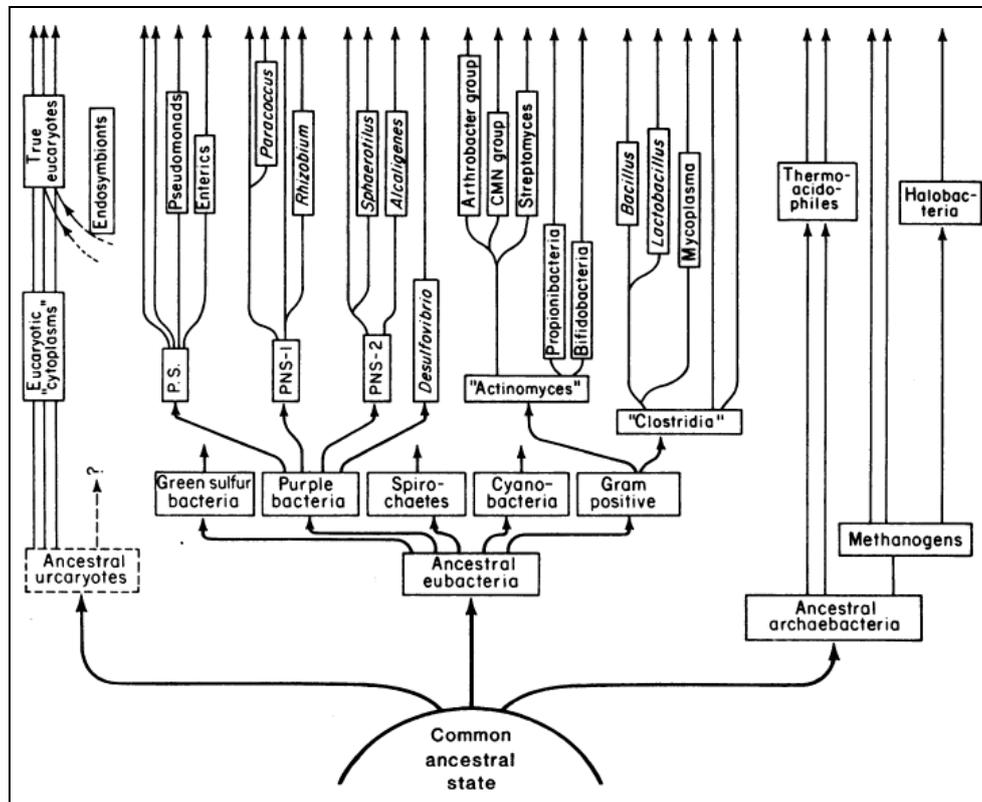


Figure 2. Représentation des grandes lignées du vivant selon Fox et collaborateurs en 1980. Les trois domaines proposés dérivent d'un même ancêtre commun (Atoji et al. 1980).

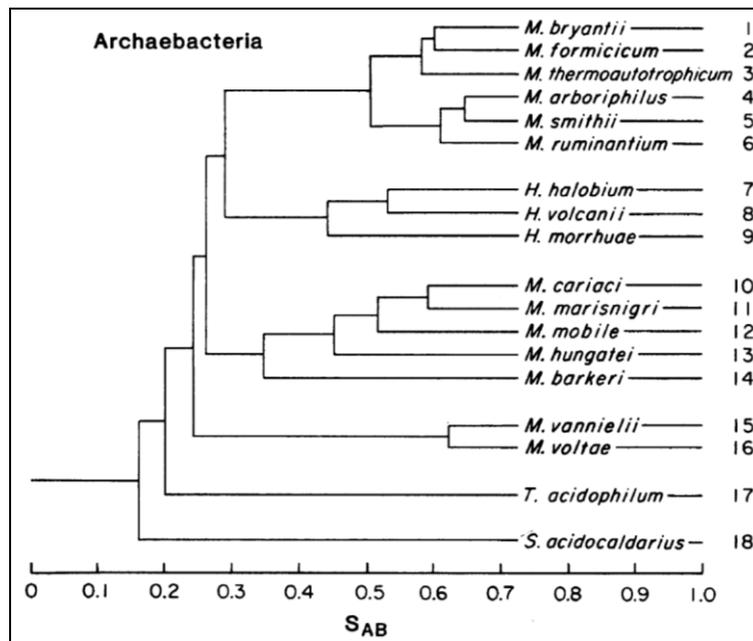


Figure 3. Premier dendrogramme représentant les relations entre les « archaeobacteria » selon Fox et collaborateurs en 1980. (Atoji et al. 1980). La racine est placée entre un groupe contenant les méthanogènes (1-6, 10-16), les halophiles (7-9) et *Thermoplasma acidophilum* (17), et les Sulfolobales, représentés par *Sulfolobus acidocaldarius* (18). Ces deux groupes correspondent aux futurs phyla des Euryarchaeota et des Crenarchaeota. Au sein du premier groupe, on voit déjà apparaître trois groupes de méthanogènes correspondant aux ordres des Methanobacteriales (1-6), Methanomicrobiales (10-14) et Methanococcales (15-16). (Atoji et al. 1980).

Le début des années 1980 a vu le remplacement des catalogues d'oligonucléotides par les premiers séquençages complets de molécules d'ARNr SSU d'« archaeobactéries » (R. Gupta, Lanter, and Woese 1983; Olsen et al. 1985). Il devient donc possible de construire des arbres phylogénétiques du vivant à partir de la séquence primaire de cette molécule (N. R. Pace, Olsen, and Woese 1986) et ainsi de confirmer la phylogénie proposée précédemment sur la base du S_{AB} . La [Figure 4](#) montre un premier arbre non raciné inféré à partir de 950 positions des séquences de l'ARNr SSU de 21 espèces, où l'on peut observer les trois domaines « Eucaryotes », « Eubacteria » et « Archaeobacteria ». L'ARNr SSU semblait être un marqueur efficace pour l'analyse de la phylogénie du vivant, et particulièrement des procaryotes. Cette molécule est présente chez tous les organismes vivants connus, sa fonction est très conservée, il est donc facile de trouver son homologue dans une nouvelle espèce, elle est assez longue pour qu'un nombre de positions suffisamment informatives puisse être utilisé et elle est peu transférée entre différents organismes. Ces caractéristiques sont aussi celles des autres ARN ribosomiques (23S et 5S chez les procaryotes), mais l'ARNr 5S, étant plus court, contient moins de signal informatif que l'ARNr SSU et un grand nombre de molécules d'ARNr SSU était déjà séquencé. L'ARNr SSU s'est de fait imposé comme marqueur pour l'étude de la phylogénie des microorganismes. La possibilité d'avoir un seul marqueur pour étudier l'évolution de l'ensemble du vivant est une avancée sans précédent,

particulièrement en ce qui concerne l'évolution des microorganismes ; celle des bactéries et archées bien sûr, mais aussi celle des organismes unicellulaires eucaryotes, considérés comme un seul groupe, par exemple par Whittaker (Whittaker 1969). Le développement de la phylogénie moléculaire a donc ouvert la porte à l'étude de l'évolution et de la phylogénie microbienne et a ainsi donné un nouveau cadre de pensée à la classification microbienne.

Il me semble intéressant d'aborder un point un peu particulier dans cette révolution portée par Woese à propos de la refonte de la classification microbienne. Il apparaît à la lecture des articles qu'il a écrit dans les années 80 que deux idées plutôt contradictoires sont abordées. D'une part, comme je viens de l'expliquer, Woese propose que la classification doive se faire uniquement sur des critères phylogénétiques, la phylogénie des microorganismes se faisant sur la base de caractères moléculaires, et particulièrement sur la séquence des ARN ribosomiques. D'autre part, la proposition du domaine des « Archaeobactéries » est suivie d'une description des caractéristiques communes entre les organismes de ce groupe. Ce qui n'est au départ (C R Woese and Fox 1977) qu'une description de ces microorganismes, devient bientôt une recherche de synapomorphies justifiant l'existence d'un groupe taxonomique à part entière. Par exemple, dans l'article « Archaeobacteria » de 1978 (C R Woese, Magrum, and Fox 1978) qui est une synthèse des travaux sur ces questions, Woese répertorie quatre grands caractères communs aux archées selon lui : les caractéristiques particulières des ARN ribosomiques et des ARN de transferts, l'absence de parois cellulaires de peptidoglycane, les lipides de membranes contenant des isoprénoïdes et reliés par des liens éthers, et l'habitat « spécialisé » des Archaea ; il se sert de ces caractéristiques pour justifier l'existence du groupe taxonomique « archaeobacteria » : « *In summary, enough comparative information now exists pertaining to the archaeobacteria to make it possible to characterize the kingdom in a positive way* » (C R Woese, Magrum, and Fox 1978). Il est facilement compréhensible qu'une proposition de changement aussi importante dans la classification microbienne ait fait débat (Wheeler, Kandler, and Woese 1992; Mayr 1990; Mayr 1991; Margulis and Guerrero 1991), et que Woese cherche à répondre à différentes critiques et à justifier la proposition du nouveau groupe « Archaeobacteria », mais il est paradoxal d'utiliser ces caractères phénotypiques à homologie discutable comme justification d'une classification qui se veut basée sur les relations phylogénétiques de caractères homologues avérés. Une autre critique peut être formulée à l'encontre de cette description, à propos des environnements « spécialisés » évoqués par Carl Woese. Certes, les archées connues à ce moment-là étaient principalement des extrémophiles, halophiles et thermoacidophiles, et des méthanogènes. Les milieux dans lesquels se développent ces organismes peuvent nous sembler particuliers, mais ils sont totalement différents entre eux et ils impliquent des adaptations différentes, par exemple, entre les organismes vivant à haute

température et ceux vivant dans des conditions de haute salinité. Leur regroupement n'est le résultat que de leur exotisme à nos yeux, ce qui constitue un argument subjectif et anthropocentré.

c. La nouvelle classification du vivant : « Archaea », « Bacteria », « Eucarya »

La vision tripartite du vivant est renforcée en 1990 quand Woese et ses collaborateurs proposent une nouvelle taxonomie en modifiant les noms des trois royaumes primaires pour les remplacer par trois domaines : « Bacteria », « Eucarya » et « Archaea » (C R Woese, Kandler, and Wheelis 1990). Le premier but de cette modification des noms est d'abord de marquer la séparation entre les Archaea et les Bacteria et ainsi de désacraliser définitivement la dichotomie Procaryotes/Eucaryotes, pour ne pas laisser de doute quant au fait que les Archaea ne sont pas une branche parmi les bactéries. L'autre but de cette proposition est d'insister sur le fait que la classification et la taxonomie doivent se faire selon un « système naturel », que le plus logique est celui de l'histoire évolutive comme l'a proposé Darwin et que, finalement, la phylogénie moléculaire est la méthode la plus efficace pour sa reconstruction : « *The time has come to bring formal taxonomy into line with the natural system emerging from molecular data* » (C R Woese, Kandler, and Wheelis 1990). Un arbre du vivant est ainsi proposé (Figure 4), inféré sur la base des séquences de l'ARNr SSU et raciné entre les bactéries d'une part et les archées et les eucaryotes d'autre part. La racine est placée ainsi d'après les analyses faites à la même époque sur des paralogues présents chez le dernier ancêtre commun, dont chaque copie peut servir de racine à la phylogénie de l'autre (Iwabe et al. 1989). La place de cette racine admet donc un ancêtre commun exclusif entre les archées et les eucaryotes et éloigne les archées des bactéries. Pour les auteurs, c'est un argument de poids pour appuyer la validité du domaine « Archaea » et le séparer du domaine « Bacteria ».

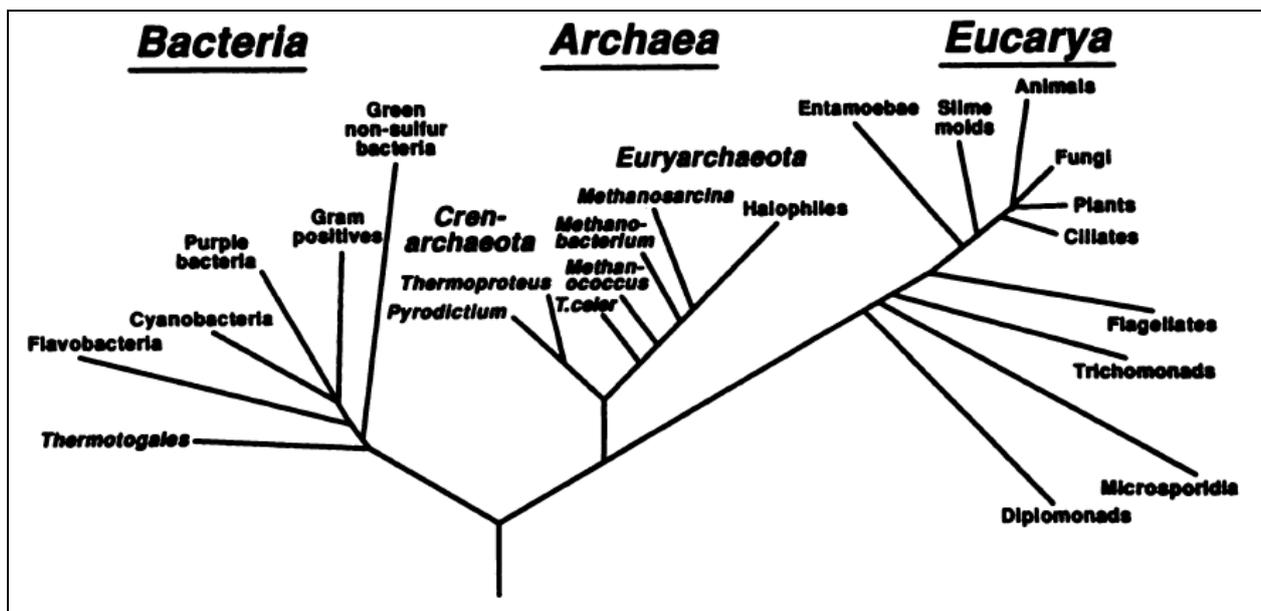


Figure 4. Arbre du vivant représentant les trois domaines Bacteria, Archaea et Eucarya. (Olsen and Woese 1993). On voit bien les deux « royaumes » archéens : Crenarchaeota et Euryarchaeota. La racine de l'arbre est placée entre les bactéries et les archées/eucaryotes d'après l'étude des phylogénies de couples de paralogues anciens.

Le domaine Archaea est composé de deux lignées majeures, appelées « royaumes » dans cette nouvelle classification, les Euryarchaeota et les Crenarchaeota, entre lesquelles se place la racine de la phylogénie des archées. Les Euryarchaeota rassemblent les halophiles extrêmes, trois lignées de méthanogènes, le genre *Archaeoglobus*, le genre *Thermoplasma*, et le groupe des *Thermococcus-Pyrococcus*, deux groupes de thermophiles. Les Euryarchaeota sont donc phénotypiquement très divers, avec des modes de vie très différents, en termes écologiques (thermophiles, acidophiles, halophiles...) et métaboliques (méthanogènes, sulfato-réducteur, hétérotrophes...), ce qui leur a valu leur nom d'Euryarchaeota, « eury » du grec « εὐρύς » pour « large ». Les Crenarchaeota regroupent uniquement des organismes thermophiles extrêmes et ont un mode de vie similaire, nécessitant du soufre comme source d'énergie. Leur nom vient de « cren » du grec « κρήνη » pour « source », rappelant leur ressemblance phénotypique avec l'ancêtre commun des archées, que les auteurs supposent être plus proche de ce groupe là. Cette classification des archées en deux royaumes Euryarchaeota et Crenarchaeota restera la norme dans ses grandes lignes pendant presque 20 ans.

2. A partir des années 1990 : Exploration de la diversité archée.

a. Etudes environnementales

Le début des années 1990 a été marqué par la découverte concomitante par DeLong (DeLong 1992) et Fuhrman et collaborateurs (Fuhrman, McCallum, and Davis 1992), d'archées marines mésophiles et aérobies, phylogénétiquement proches pour certaines des Crenarchaeota (le groupe I), et pour d'autres se plaçant au sein des Euryarchaeota (le groupe II) (Figure 5).

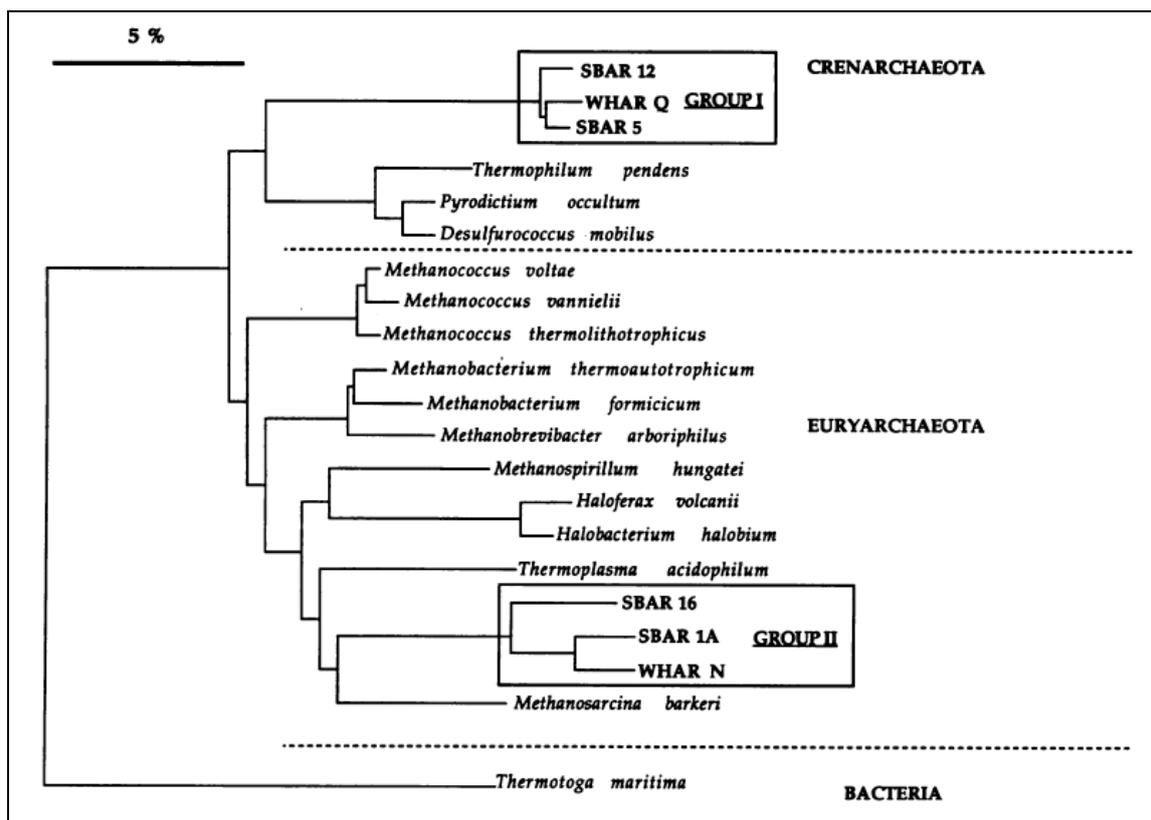


Figure 5. Phylogénie des Archées montrant les groupes I et II. Arbre inféré à partir de l'ARNr SSU d'archées cultivées et de séquences d'espèces non-cultivées d'échantillons marins (DeLong 1992). Le groupe I se place à la base des Crenarchaeota déjà connues, et le groupe II se place à l'intérieur des Euryarchaeota, comme groupe frère des Methanosarcinales. Cet arbre est raciné sur la bactérie *Thermotoga maritima*.

Cette découverte marque le début des études environnementales pour l'exploration de la diversité des archées. En effet, en ayant mis en évidence que l'ARNr SSU était un bon marqueur pour la phylogénie des organismes, Woese a ouvert la porte à l'étude de la diversité microbienne grâce à une approche moléculaire par amplification par PCR et séquençage des gènes d'ARNr SSU à partir de l'ADN d'échantillons naturels (Olsen et al. 1986).

DeLong, Fuhrman et leurs collaborateurs ont extrait l'ADN de différents échantillons marins de surface, puis amplifié et séquençé l'ARNr SSU et ont ainsi pu avoir un aperçu de la diversité

microbienne présente. Les archées ayant été décrites jusqu'alors comme un domaine d'organismes vivant dans des environnements « spécialisés », elles n'étaient pas attendues dans des échantillons d'eau de mer de surface, oxygénée et plutôt froide, et la première hypothèse avancée par les auteurs de ces deux articles pour expliquer la présence de séquences d'archées dans leurs échantillons est une possible contamination. Il se trouve que des méthanogènes avaient été découverts peu de temps avant dans l'intestin de gros animaux, dont l'homme, et qu'une contamination humaine aurait pu expliquer la présence des archées du groupe II (se plaçant au sein des Euryarchaeota dont font partie les méthanogènes). Cependant, les auteurs vont privilégier une autre explication, celle d'une diversité d'archées vraiment inféodées au milieu planctonique marin, malgré le caractère exceptionnel de cette découverte comme le montre cette phrase de DeLong : « *A more provocative explanation is that group I and group II represent undescribed mesophilic, aerobic members of the archaea* » (DeLong 1992). En effet, ces deux travaux décrivant une même diversité sur des échantillons océaniques provenant de quatre localités géographiques différentes sur les deux côtes des Etats-Unis (DeLong 1992; Fuhrman, McCallum, and Davis 1992) se confirment l'un l'autre et font apparaître une diversité d'archées inconnues et non cultivées assez remarquable. Ces communautés d'archées semblent alors importantes en termes de biomasse dans les océans, bien que moins importantes que les bactéries présentes dans les mêmes milieux. Cette découverte est très importante d'un point de vue écologique puisqu'elle met en évidence une part de la biodiversité océanique ignorée jusqu'alors. Elle est aussi importante du point de vue de la connaissance de la diversité des archées, car elle décrit des groupes totalement inconnus dont le mode de vie mésophile et aérobie défie l'idée selon laquelle les archées vivent exclusivement dans des environnements « spécialisés ». Enfin, du point de vue de leur histoire évolutive, la découverte selon laquelle l'une de ces lignées se place à la base des Crenarchaeota ouvre de nombreuses questions quant à la nature de l'ancêtre commun des archées. Le groupe I sera d'ailleurs proposé comme étant un phylum à part entière en 2008 par Brochier-Armanet et collaborateurs sous le nom de « Thaumarchaeota », du grec « *θαυμασ* » « thaumas » pour « merveilleux » (Brochier-Armanet et al. 2008).

Les études environnementales sur la diversité des archées vont mettre en évidence dans les années suivantes un nombre important de lignées non cultivées vivant, pour certaines, dans des environnements nous semblant moins extrêmes que ceux des archées cultivées et connues (Christa Schleper, Jurgens, and Jonuscheit 2005). À l'instar de l'océan, on trouve des archées dans des sédiments marins, des lacs d'eau douce, des sols, en symbiose avec d'autres organismes, mais aussi dans des milieux très acides (mines acides), ou des tapis microbiens riches en méthane, en association avec des méthanogènes et des bactéries sulfato-réductrices. Initialement basées sur l'amplification et le séquençage de l'ARNr SSU, les études environnementales se sont

progressivement diversifiées en termes de techniques. Le groupe de DeLong a largement contribué à ce développement, par exemple avec l'isolement et l'analyse d'un large fragment génomique de 40 kilobases d'une archée du groupe I, un des premiers exemples de l'application des techniques de métagénomique (Stein et al. 1996). Dans ce travail, les auteurs ont extrait et digéré l'ADN d'un échantillon de plancton marin pour construire une banque de BAC (Bacterial Artificial Chromosome) modifiés pour recevoir des gros inserts d'ADN. Ces inserts pouvant contenir plusieurs gènes, si un marqueur phylogénétique tel que l'ARNr SSU est présent, il est possible, non seulement de savoir à quel groupe taxonomique appartient l'organisme, mais surtout d'étudier les gènes placés autour de ce marqueur et ainsi d'en connaître un peu plus sur les activités métaboliques que cet organisme peut réaliser. L'application de cette méthodologie à des environnements dans lesquels une seule espèce d'archée est présente a permis de reconstruire et de séquencer des génomes presque complets, comme ça a été le cas pour l'organisme *Cenarchaeum symbiosum* en 2006, vivant en symbiose avec une éponge et une bactérie (Hallam et al. 2006). Les auteurs de ce papier (dont DeLong) ont pour cela construit une banque de fosmides, séquencé chaque insert et réassemblé le génome en utilisant les parties chevauchantes des différents fragments. Ce génome a été le premier à être séquencé pour une archée du groupe I.

Il est également possible d'étudier la grande diversité d'archées non cultivées en ayant recours aux observations par FISH (hybridation in situ en fluorescence), en ciblant spécifiquement les ARNr SSU d'archées (Christa Schleper, Jurgens, and Jonuscheit 2005; DeLong 1992). Cette technique a pour avantage de montrer l'abondance de la population ciblée, mais aussi de montrer des cellules vivantes et de prouver que la détection d'ARNr SSU d'archées dans ces milieux n'est pas un artefact lié à des contaminations ou des cellules mortes (Figure 6). C'est ainsi que Karner, DeLong et Karl ont pu montrer en 2001 que les crenarchaeota marines du groupe I étaient très abondantes dans les milieux océaniques, de l'ordre de 20% de la population microbienne planctonique (Karner, DeLong, and Karl 2001).

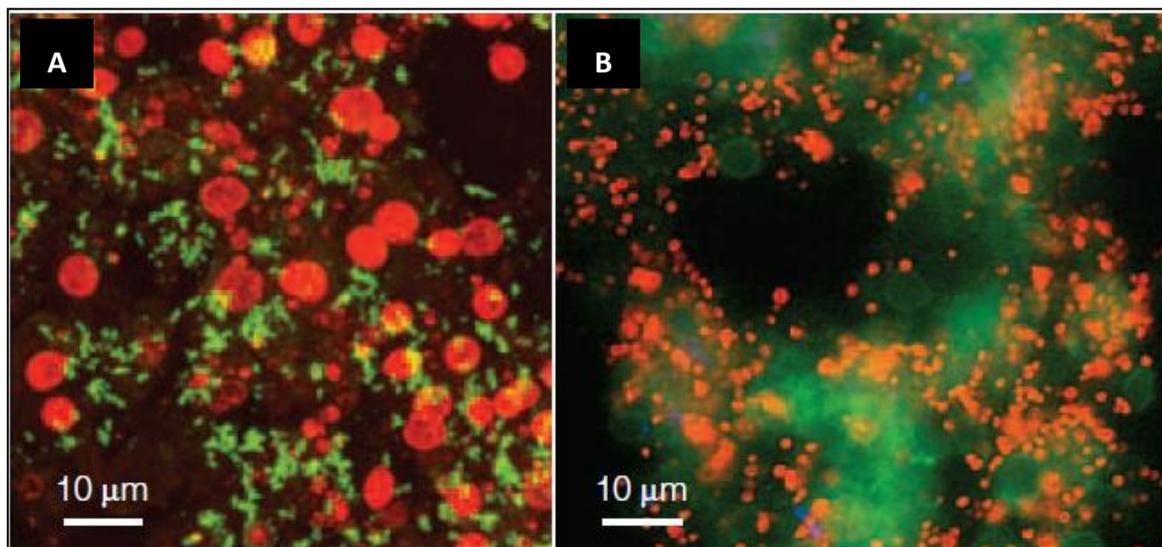


Figure 6. Images d'archées non cultivées en FISH. Issues de l'article Schleper et collaborateurs en 2005 (Christa Schleper, Jurgens, and Jonuscheit 2005). A : cellules de *Cenarchaeum symbiosum*, archée du groupe I (en vert), et d'*Axinella mexicana*, une éponge, hôte de *C. symbiosum* (en rouge). Image de C. Preston, Monterey Bay Aquarium Research Institute, Moss Landing, USA. B : Enrichissement de crenarchaeota du groupe I (en rouge), sur des racines de tomates, et des bactéries associées (en vert). Image de H. Simon, Oregon Health & Science University, Beaverton, USA.

Si la découverte d'archées mésophiles et aérobies fut très importante pour la perception que l'on pouvait avoir de ce domaine, la recherche d'archées dans des milieux extrêmes a aussi apporté son lot de lignées inconnues, certaines se plaçant à proximité de la base de la phylogénie des archées. En 1994, Barns et ses collaborateurs (Barns, Fundyga, Jeffries, & Pace, 1994; Barns, Delwiche, Palmer, & Pace, 1996) ont décrit la diversité d'archées dans l'eau et les sédiments de deux sources chaudes du parc de Yellowstone. Différentes séquences de Crenarchaeota ont été décrites, branchant à proximité des Desulfurococcales et des Sulfolobales ou au sein des Thermoproteales. D'autres séquences branchaient avec le groupe I nouvellement décrit par DeLong et Furhmann, et d'autres encore branchaient au sein de l'arbre des archées avant même la divergence entre Crenarchaeota et Euryarchaeota. Cette position a amené les auteurs à proposer que ces organismes pouvaient être les représentants d'un nouveau phylum d'archées et ils proposèrent le nom de « Korarchaeota », du grec « κορος » pour « jeune homme », pour la divergence précoce de ce groupe dans l'histoire des archées (Barns et al. 1996).

Au cours des années 2000, un nouveau « type » d'archées a été découvert : de très petites archées, à commencer par *Nanoarchaeum equitans* en 2002 (H. Huber et al. 2002). Cette archée de seulement 400 nm de diamètre vit en symbiose avec une autre archée, *Ignicoccus hospitalis*, une crenarchée. Son génome ne fait que 0,5 megabases, et son séquençage fera apparaître de nombreux traits particuliers tels que l'absence d'un grand nombre de gènes essentiels chez les archées ou des gènes coupés en deux parties, notamment des gènes codants des ARNt. Une petite taille cellulaire et

un génome de taille très réduite sont des caractéristiques classiques rencontrées chez les bactéries endosymbiotiques ou parasites, mais n'avaient jamais été observées chez les archées ou chez des hyperthermophiles. Les analyses phylogénétiques de l'ANRr SSU ont placé *N. equitans* très basalement dans l'arbre des archées mais avec un faible soutien (H. Huber et al. 2002). Ceci a conduit les auteurs à proposer qu'elle serait la représentante d'un nouveau phylum d'archées, les « Nanoarchaeota » (« nano » pour sa petite taille cellulaire de l'ordre du nanomètre). Plus récemment, d'autres archées de petite taille ont été décrites, le groupe des ARMAN en 2006 par Baker et collaborateurs (Baker et al. 2006) et les Nanohaloarchaeota en 2011 par Narasingarao et collaborateurs (Narasingarao et al. 2012). Le groupe des ARMAN pour « archaeal Richmond Mine acidophilic nanoorganisms » a été trouvé dans les mines acides de Richmond, au niveau de biofilms principalement dominés par des archées de type thermoplasmatales (*Thermoplasma* et *Ferroplasma*) et des bactéries du genre *Leptospirillum*. En 2010, trois génomes ont été publiés, reconstruits à partir de données métagénomiques (Baker et al. 2010) : '*Candidatus* Micrarchaeum acidiphilum ARMAN-2', '*Candidatus* Parvarchaeum acidiphilum ARMAN-4' et '*Candidatus* Parvarchaeum acidophilus ARMAN-5'. Ces deux derniers génomes sont très proches en termes de similarité de séquences. Dans les trois cas, il s'agit de toutes petites archées comparables à *N. equitans* en termes de taille (environ 500 nm de diamètre), mais leurs génomes sont deux fois plus gros (environ 1 megabase) bien qu'ils aient des caractéristiques comparables telles que l'absence de nombreux gènes essentiels ou des gènes splittés chez '*Ca.* Micrarchaeum acidiphilum ARMAN-2'. Les absences de gènes demanderont à être confirmées car ces génomes ne sont pas totalement complets. Leurs caractéristiques cellulaires et génomiques portent à croire qu'elles ont un type de vie parasitique ou symbiotique mais n'étant pas cultivées, ce point reste sans réponse. Enfin, les Nanohaloarchaea sont de petites archées halophiles dont les génomes ont été séquencés d'après des données métagénomiques. C'est le cas pour '*Candidatus* Nanosalinarum sp.' et '*Candidatus* Nanosalina sp.' (Narasingarao et al. 2012) dont le séquençage a été suivi de près par celui de '*Candidatus* Haloredivivus' (Ghai et al. 2011). Ces trois génomes ont une taille d'environ 1,2 megabases et '*Ca.* Nanosalinarum sp.' et '*Ca.* Nanosalina sp.' ont des cellules de très petite taille, environ 600 nm. Les analyses phylogénétiques ont placé ces organismes au sein des Euryarchaeota, comme groupe frère des Halobacteria. C'est pourquoi la nouvelle classe des Nanohaloarchaea a été proposée par Narasingarao et collaborateurs (Narasingarao et al. 2012).

L'un des derniers phyla d'archée à avoir été proposé est celui des Aigarchaeota par Nunoura et collaborateurs en 2011 (Nunoura et al. 2011), lors de la publication du génome composite de '*Candidatus* Caldiarchaeum subterraneum'. Cet organisme appartient au groupe HWCGI (Hot Water Crenarchaeotic Group I), un groupe d'organismes thermophiles, non cultivés, et

phylogénétiquement proche des Thaumarchaeota et des Crenarchaeota tout en étant distinct. Cette position basale dans l'arbre des archées et des caractéristiques génomiques particulières, telles que la présence de protéines typiquement eucaryotes intervenant dans le système ubiquitine-protéasome (système de marquage des protéines par l'ubiquitine en vue de leur dégradation par le protéasome) ont conduit les auteurs à proposer la création du nouveau phylum. Ils proposent le nom d'« Aigarchaeota », du grec « αἴγι », « aigi », pour « aurore », ces organismes présentant des caractéristiques intermédiaires entre les hyperthermophiles et les mésophiles dans cette partie de l'arbre des archées.

Depuis, le phylum « Geoarchaeota » ou « NAG1 » (pour « novel archaeal group 1 ») a été proposé par Kozubal et collaborateurs en 2013 (Kozubal et al. 2013), d'après l'analyse de quatre génomes construits à partir de métagénomes issus de tapis microbiens du parc de Yellowstone. Cet environnement est acide, a une haute température (65-80°C), contient du fer ferreux et une faible concentration d'oxygène. Ces organismes se placeraient en groupe frère des Crenarchaeota et auraient des caractéristiques génomiques particulières par rapport aux autres archées, ce qui a conduit les auteurs à proposer la création d'un nouveau phylum. Néanmoins, leur position phylogénétique demanderait à être confirmée, particulièrement au regard de leur position basale supposée dans l'arbre des archées.

b. Les années 2000 : Séquençage de nombreux génomes complets

Depuis la découverte des archées en tant que domaine à part entière, la phylogénie moléculaire a toujours été un outil privilégié pour leur étude. Le séquençage de génomes complets a donc représenté une grande avancée, de façon générale en biologie mais particulièrement pour la connaissance de ce domaine. En effet, peu d'espèces d'archées sont cultivées par rapport aux eucaryotes et aux bactéries, et peu d'outils de génétique et de biologie moléculaire ont pu être développés du fait de leurs particularités cellulaires, de conditions de vie et même génétiques (bien que de plus en plus de modèles aient été développés ces dernières années (Leigh et al. 2011)). L'accès aux séquences complètes de leurs génomes a apporté une quantité de données sans précédent sur ces organismes encore mal connus. *Methanococcus jannaschii* fut la première archée dont le génome a été complètement séquencé en 1996 (Bult et al. 1996). Depuis, 227 génomes d'archées ont été complètement séquencés (Kyrpides 1999)(aout 2013). Bien que ne couvrant pas la totalité de la diversité des archées, au moins un génome représentant chaque phyla est disponible à l'heure actuelle (Euryarchaeota, Crenarchaeota, Korarchaeota, Thaumarchaeota, et Nanoarchaeota). Pour certaines espèces, plusieurs souches sont séquencées ce qui permet des études de génomique

comparée à petite échelle taxonomique. L'exemple le plus important est celui de l'espèce *Sulfolobus islandicus* dont 15 souches sont séquencées. Certaines ont été étudiées par génomique comparative par Reno et collaborateurs en 2009 (Reno et al. 2009) d'un point de vue biogéographique.

La publication de génomes complets d'archées et l'augmentation en nombre et en diversité est aussi très importante en ce qui concerne l'étude de leur histoire évolutive. Jusqu'ici, le marqueur phare était l'ARNr SSU, une molécule privilégiée pour la reconstruction phylogénétique de l'histoire des organismes, très conservé au niveau de sa séquence et de sa fonction, et peu transféré. Il reste néanmoins un marqueur unique et à ce titre soumis à certains biais tels qu'un nombre limité de positions, donc de caractères, pouvant de plus accumuler une saturation mutationnelle importante. L'accès aux séquences de gènes codant des protéines grâce aux génomes complets a donc permis d'utiliser de nouveaux marqueurs pour reconstruire la phylogénie des archées. Les protéines ribosomiques ont été parmi les premières séquences utilisées, puisque, intervenant dans le même système biologique que l'ARNr SSU, elles sont aussi très conservées et peu transférées, ce qui en fait des marqueurs phylogénétiques de choix. C'est d'ailleurs aussi le cas pour les sous-unités de l'ARN polymérase. Et plus que l'analyse de nouveaux marqueurs, l'analyse conjointe de plusieurs protéines a été rendue possible par l'arrivée des données génomiques. Ces études ont été commencées par Matte-Tailliez et collaborateurs en 2002 (Matte-Tailliez et al. 2002), et améliorées depuis par l'ajout de nouvelles protéines et l'augmentation de l'échantillonnage taxonomique (Simonetta Gribaldo and Brochier-Armanet 2006). Ces travaux ont permis de confirmer les grandes lignes de la phylogénie obtenue par Woese en 1980 et de raffiner la phylogénie interne des différents phyla d'archées.

Les génomes complets ont tout d'abord été obtenus à partir d'organismes cultivés pour des raisons évidentes de quantité de matériel génétique nécessaire au séquençage et de sécurité quant aux contaminations. Ainsi que je l'ai déjà évoqué, des génomes peuvent être séquencés pour des organismes non isolés en culture pure à partir de données métagénomiques, comme par exemple *Cenarchaeum symbiosum* (Hallam et al. 2006). De plus en plus de génomes sont reconstruits de cette manière, par exemple ceux des trois archées appartenant au groupe ARMAN (Baker et al. 2010), ou à partir d'une culture enrichie pour celui de '*Candidatus* Korarchaeum cryptofilum' (Elkins et al. 2008). Cette stratégie mène normalement à l'obtention de ce qu'on appelle un génome composite. Aujourd'hui, le séquençage à partir de cellules uniques est devenu possible et le premier génome ainsi séquencé, celui de la bactérie '*Candidatus* Sulcia muelleri DMIN', a été publié en 2010 (Woyke et al. 2010). La première archée séquencée de cette façon a été '*Candidatus* Nitrosoarchaeum limnia SFB1' en 2011 (Blainey et al. 2011), un membre des Thaumarchaeota (anciennement crenarchaeota du groupe I), dont les représentants sont très difficiles à cultiver

encore aujourd'hui. Cette nouvelle technique ouvre des portes quant à la connaissance de l'ensemble des groupes non cultivés d'archées, comme le montre la publication très récente de Rinke et collaborateurs, décrivant le séquençage de 201 archées et bactéries, principalement issues de groupes peu connus (Rinke et al. 2013). L'accès à leurs séquences génomiques pourrait permettre d'en savoir plus sur leur évolution, bien sûr, mais aussi sur leurs métabolismes potentiels et, peut être, sur les conditions nécessaires à leur développement, ce qui pourrait servir à mettre en place des milieux propices à leur culture pour une étude *in vivo* plus approfondie.

3. Conclusion

La découverte des Archaea en tant que domaine du vivant à part entière est étroitement liée à l'étude de l'évolution grâce à des données moléculaires, très vite remplacées par la phylogénie moléculaire. Les difficultés de culture d'une grande partie de leur diversité ont aussi renforcé l'importance des données génomiques dans l'étude des archées. Ces raisons historiques expliquent que l'étude des archées soit intimement liée à celle de leur histoire évolutive et de leurs génomes ; c'est pourquoi il m'a parut important de rappeler les grandes étapes historiques qui ont conduit à la connaissance actuelle de ce domaine du vivant, et l'importance de l'étude de leur évolution.

B. Phylogénie des Archaea

1. Diversité des Archaea : la phylogénie de l'ARNr SSU

Comme nous venons de le voir six phyla d'archées sont actuellement relativement acceptés : les Euryarchaeota et les Crenarchaeota, les deux phyla originellement proposé par Woese en 1980, les Korarchaeota, les Nanoarchaeota, les Thaumarchaeota et les Aigarchaeota. Ces six phyla regroupent des organismes très divers et j'aimerais ici présenter une vue d'ensemble non exhaustive de cette diversité. La Figure 7, ci-dessous, est un arbre publié par Schleper et collaborateurs en 2005 (Christa Schleper, Jurgens, and Jonuscheit 2005), inféré sur les séquences d'ARNr SSU d'archées cultivées et non cultivées. Bien que datant de 2005, cette phylogénie est très complète en ce qui concerne les groupes non cultivés et donne un bon aperçu de ce que l'on connaît de la diversité des archées actuellement.

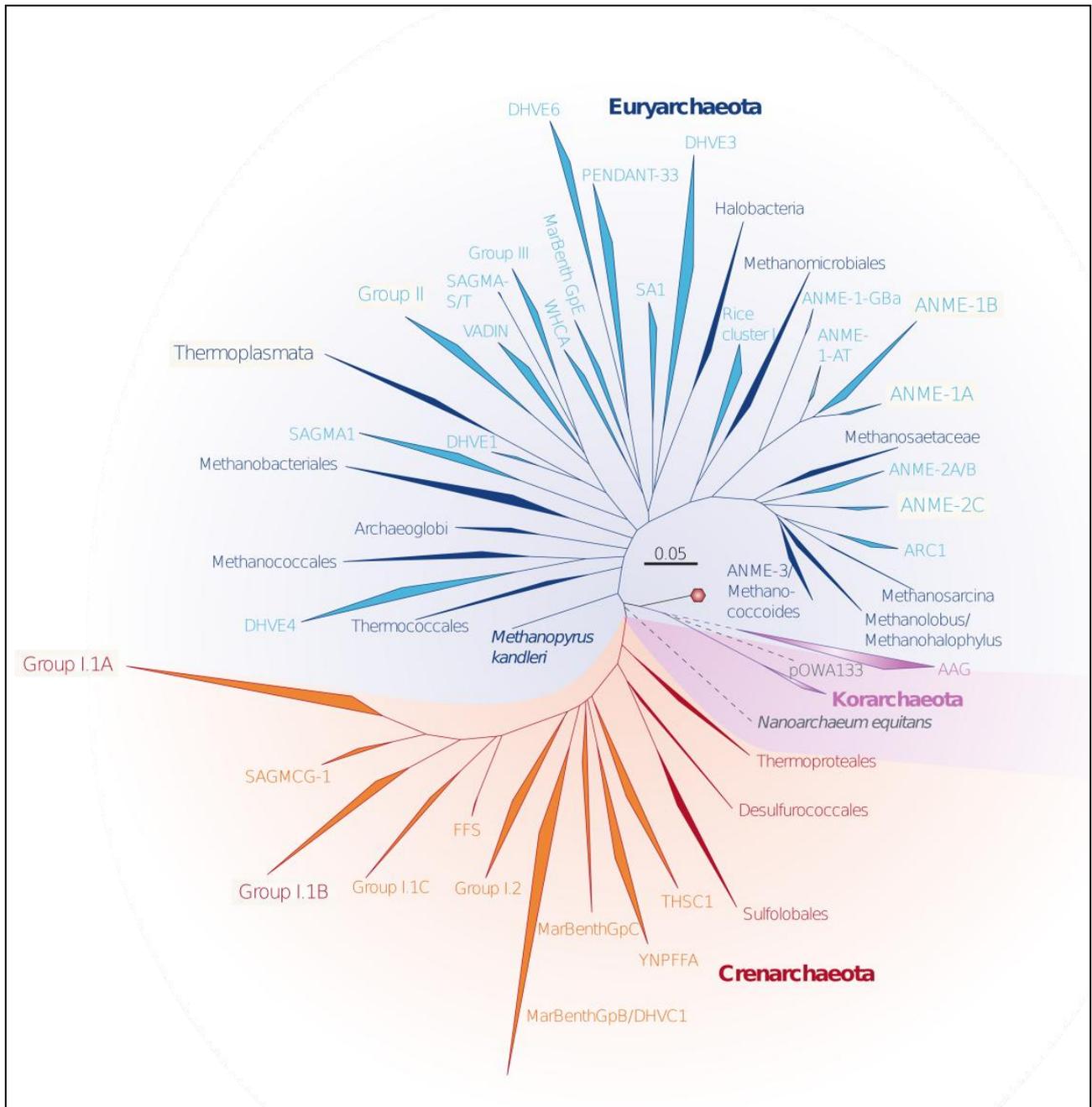


Figure 7. Phylogénie des archées cultivées et non cultivées. (Christa Schleper, Jurgens, and Jonuscheit 2005). Inféré à partir de l'ARNr SSU, cet arbre représente la diversité du domaine Archaea.

Dans cette phylogénie, le phylum Thaumarchaeota n'apparaît pas en tant que tel mais regroupe le groupe I et d'autres branches comme les SAGMCG-1. Les séquences du phylum Aigarchaeota se placeront groupe frère de ces groupes et juste après les séquences de crenarchées lors de leur découverte (Nunoura et al. 2010). Par contre, on voit bien apparaître les autres phyla, Euryarchaeota, Nanoarchaeota (placé sur une branche en pointillés car sa place n'est pas résolue), Korarchaeota et Crenarchaeota. Le groupe AAG « Ancient Archaeal Group » et le clone pOWA133 (Takai and Horikoshi 1999) représentent des séquences environnementales isolées dans une source

hydrothermale sous-marine profonde au Japon. Les analyses phylogénétiques basées sur l'ARNr SSU les placent en groupe frère d'un groupe contenant les Crenarchaeota, Korarchaeota et Thaumarchaeota (Takai and Horikoshi 1999). Des études supplémentaires seraient nécessaires pour évaluer l'importance de ce groupe, malheureusement très peu d'informations sont disponibles.

a. Les Euryarchaeota

Cette phylogénie d'ARNr SSU ([Figure 7](#)) montre neuf ordres d'Euryarchaeota (en bleu foncé) pour lesquels des génomes complets sont séquencés et des organismes cultivés : les *Methanopyrales* (représentés par *Methanopyrus kandleri*), les *Thermococcales*, les *Archaeoglobales*, les *Methanobacteriales*, les *Thermoplasmatales*, les *Halobacteriales*, les *Methanomicrobiales*, les *Methanosarcinales* et les *Methanococcales*. Depuis ont été ajoutés les ordres des *Methalocellales* (représentés ici par le groupe « Rice cluster I ») (Sakai et al. 2008) et des *Methanoplasmatales* (Mihajlovski, Alric, and Brugère 2008; Paul et al. 2012) et la classe des *Nanohaloarchaea* (Narasingarao et al. 2012). Le groupe des ARMAN a été défini sur la base de séquences environnementales et depuis, trois génomes ont été publiés (Baker et al. 2010), mais il n'est pas défini formellement comme un ordre taxonomique. Un organisme du groupe DHVE2 a été cultivé et séquencé, *Aciduliprofundum boonei* (Reysenbach et al. 2006) et un génome du groupe II d'archées marines a été reconstruit à partir de données métagénomiques (Iverson et al. 2012). À ces groupes se rajoute un nombre important de groupes d'organismes non cultivés, découverts lors d'études environnementales telles que celles décrites plus tôt. On trouve entre autre le groupe III des archées marines, découvertes par DeLong et Furhman et collaborateurs en 1992, les groupes ANME (pour Anaerobic Methanotrophic archaea) d'archées méthanotrophes, phylogénétiquement proches des archées méthanogènes, ou différents groupes DHVE, pour « Deep-sea Hydrothermal Vent Euryarchaeotic ».

Hyperthermophiles

En termes d'environnement, on trouve des Euryarchaeota dans quasiment tous les milieux. Puisque la réputation des archées est fondée sur leur extrêmophilie, commençons par les milieux hyperthermophiles ([Figure 8](#)). Différentes classes d'euryarchées vivent dans de tels milieux, i.e. à des températures supérieures ou égales à 80°C ([Figure 9](#)). Les *Thermococcales* sont tous hyperthermophiles, et ont été observés dans des sources hydrothermales terrestres ou sur des cheminées hydrothermales sous-marines, de surface ou profondes. Une espèce a même été isolée depuis un puits de pétrole, *Thermococcus sibiricus* (Mardanov et al. 2009). Leurs températures de

croissance optimale varient entre 80 et 100°C. Ce sont des organismes anaérobies, chimiorganotrophes ou chimiolithotrophes, utilisant le soufre comme accepteur final d'électrons. Les *Archaeoglobales* sont aussi des organismes hyperthermophiles anaérobies vivant dans des sédiments de sources chaudes ou des cheminées hydrothermales. Les espèces du genre *Archaeoglobus* sont sulfato-reductrices, propriété relativement rare dans le vivant, partagée seulement avec quelques bactéries comme les *Desulfurovibrio* (von Jan et al. 2010). L'organisme *Ferroglobus placidus*, archaeoglobale isolé dans un système hydrothermal sous-marin peu profond à Volcano en Italie, n'est pas sulfato-réducteur, mais est capable d'oxyder le fer à pH neutre alors que cette capacité n'avait été observée que dans des milieux soit à pH acide, soit à basse température (Hafenbradl et al. 1996). Des groupes d'organismes non cultivés tels que les groupes DHVE, ont été détectés sur des fumeurs sous-marins et un organisme du groupe DHVE2 a été isolé et séquencé, *Aciduliprofundum boonei*, phylogénétiquement proche des Thermoplasmatales et thermoacidophile (Reysenbach et al. 2006).

Un certain nombre de méthanogènes sont aussi hyperthermophiles, comme certains *Methanococcales*, *Methanobacteriales* et *Methanopyrus kandleri* (Kurr et al. 1991), seul représentant des *Methanopyrales* dont le génome ait été séquencé et pouvant vivre jusqu'à une température de 110°C.



Figure 8. Champ hydrothermal d'El Tatio au Chili. Photographies de Purificación López-García.

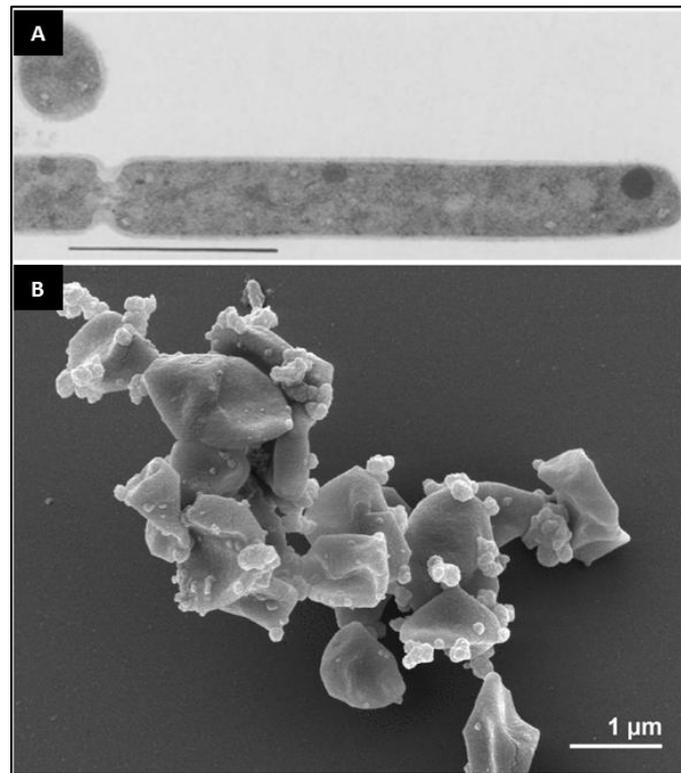


Figure 9. Euryarchées hyperthermophiles. A : Coupe fine de *Methanopyrus kandleri* (Kurr et al. 1991) Barre=1 μ m. B : Cellules d'*Archaeoglobus profundus*, image de microscopie électronique à balayage (von Jan et al. 2010).

Mésophiles et psychrophiles

A contrario, on trouve aussi des euryarchées dans des milieux extrêmement froids. Les groupes d'archées non cultivées II et III découverts dans les océans (DeLong 1992; Fuhrman, McCallum, and Davis 1992) ont été la première preuve que des archées sont présentes dans ces milieux océaniques froids, l'eau pouvant avoir une température de l'ordre de 4°C. En réalité, ces milieux froids ou tempérés (jusqu'à 50°C) sont plutôt majoritaires sur terre, et les organismes mésophiles (vivant entre 4 et 50°C) et psychrophiles (vivant à des températures inférieures à 4°C) sont très abondants, y compris parmi les archées (Cavicchioli 2006). L'espèce halophile *Halorubrum lacusprofundi* a été la première archée détectée dans ce genre de milieux extrêmement froid, en l'occurrence à *Deep Lake* en Antarctique (P.D. Franzmann et al. 1988). Des méthanogènes ont été isolés dans des eaux douces, par exemple à *Ace Lake* en Antarctique (P D Franzmann et al. 1997; P.D. Franzmann et al. 1992), ou dans de l'eau de mer, à *Skan Bay* en Alaska (Cavicchioli 2006). Mais la majeure partie de la diversité d'Euryarchaeota vivant à basse température est encore peu connue et non cultivée.

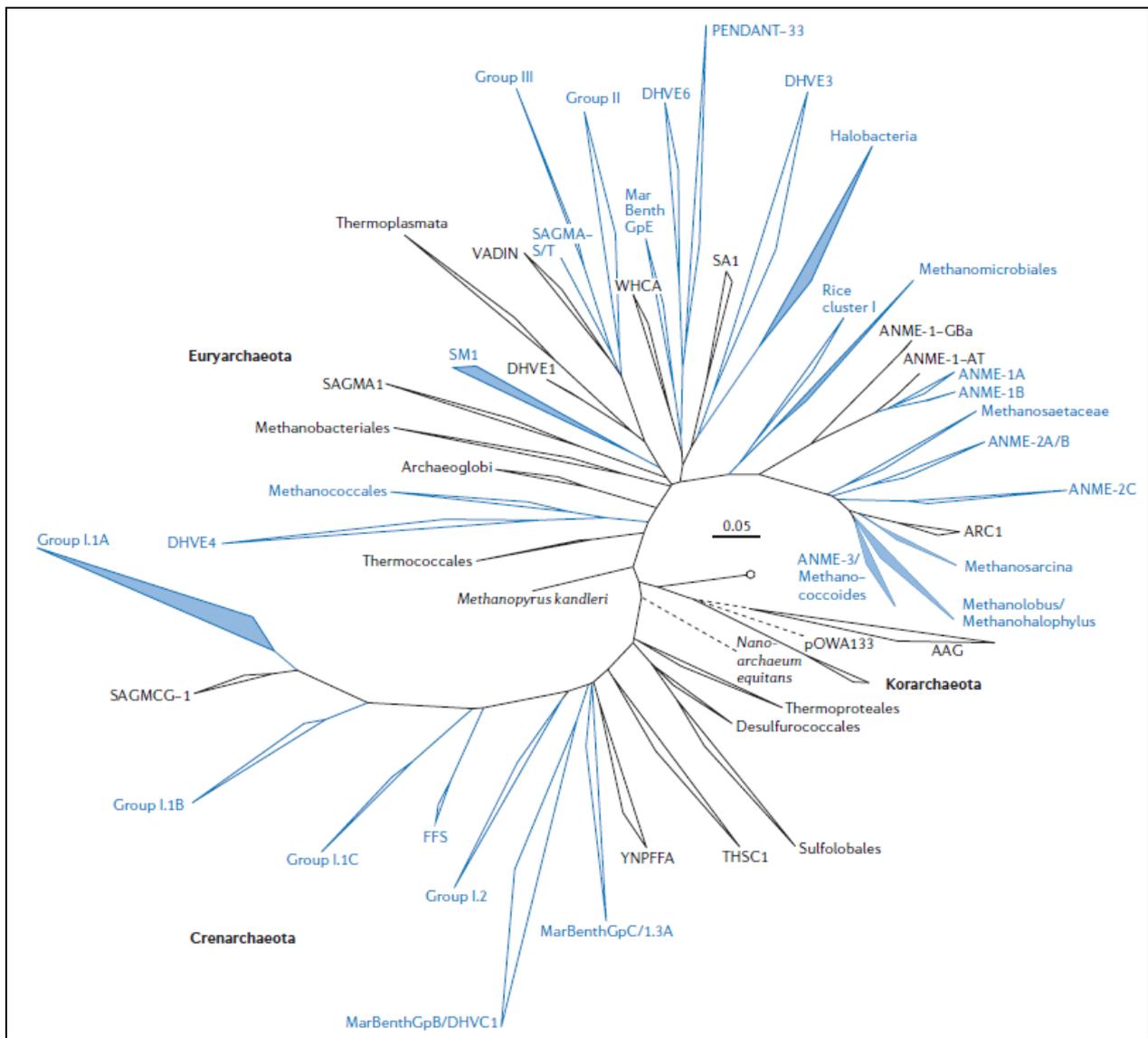


Figure 10. Diversité des archées mésophiles. Adapté de Schleper et collaborateurs 2005 (Christa Schleper, Jurgens, and Jonuscheit 2005) par Cavichioli en 2006 (Cavicchioli 2006), cet arbre montre en bleu les groupes d'archées connues pour vivre dans des environnements froids ou dont la description contient le mot « cold ».

La Figure 10 montre, en bleu, les groupes d'archées décrites dans des milieux froids. Parmi celles-ci on retrouve les archées du groupe II décrites dans des échantillons récoltés en Antarctique par DeLong en 1994 (DeLong et al. 1994), ou les archées du groupe SM1, identifiées dans des marais sulfureux à une température d'environ 10°C (Rudolph, Wanner, and Huber 2001). Sans avoir besoin de descendre à des températures aussi basses, de nombreuses archées sont observées dans des milieux mésophiles, tels que notre bouche ou des rizières pour certains méthanogènes (Knittel and Boetius 2009). La plupart des halophiles sont aussi mésophiles, et vivent par exemple dans des lacs salés ou sur la surface d'aliments soumis à salaison.

Halophiles

Un autre exemple d'extrémophilie est l'halophilie, l'adaptation à des milieux hypersalins, tels que la Mer Morte ou des aliments salés. Deux ordres d'archées phylogénétiquement proches sont particulièrement adaptés à ces milieux, les *Halobacteriales* et les *Nanohaloarchaea*. Ces organismes sont capables de croître à des concentrations en sel entre 2,5 et 5,2 M de NaCl pour les plus extrêmes, et entre 0,5 et 2,5 M pour les halophiles modérés. Les environnements dans lesquels ils se développent sont très divers, dans la mesure où la concentration en sel est forte, par exemple, *Haloterrigena turkmenica* a été isolé dans des sols (Zvyagintseva and Tarasov 1988), *Halalkalicoccus jeotgali* sur des crustacés fermentés au sel (Roh et al. 2010) et *Haloferax volcanii* vit dans la Mer Morte (Mullakhanbhai and Larsen 1975).

Les *Halobacteriales* sont pour la plupart aérobies et contiennent des pigments de type caroténoïdes, donc rouges, ce qui donne à leurs colonies des couleurs vives et caractéristiques, comme l'on peut le voir sur la [Figure 11](#). Certains sont anaérobies et photohétérotrophes grâce à la bactériorhodopsine, une pompe à proton spécifique à ce groupe, utilisant la lumière comme source d'énergie. Le gradient de protons généré permet la synthèse d'ATP (Hartmann, Sickinger, and Oesterhelt 1980). Une caractéristique génomique de ce groupe est la présence de plusieurs chromosomes ou de nombreux megaplasmides pouvant contenir jusqu'à plusieurs centaines de gènes, alors que le génome de la plupart des archées se trouve sur un seul chromosome et, potentiellement, quelques plasmides de petite taille. Enfin, on peut rappeler ici que la première archée décrite en 1880 par Farlow était une halophile, *Halococcus* (Farlow 1880; KOCUR and Hodgkiss 1973).

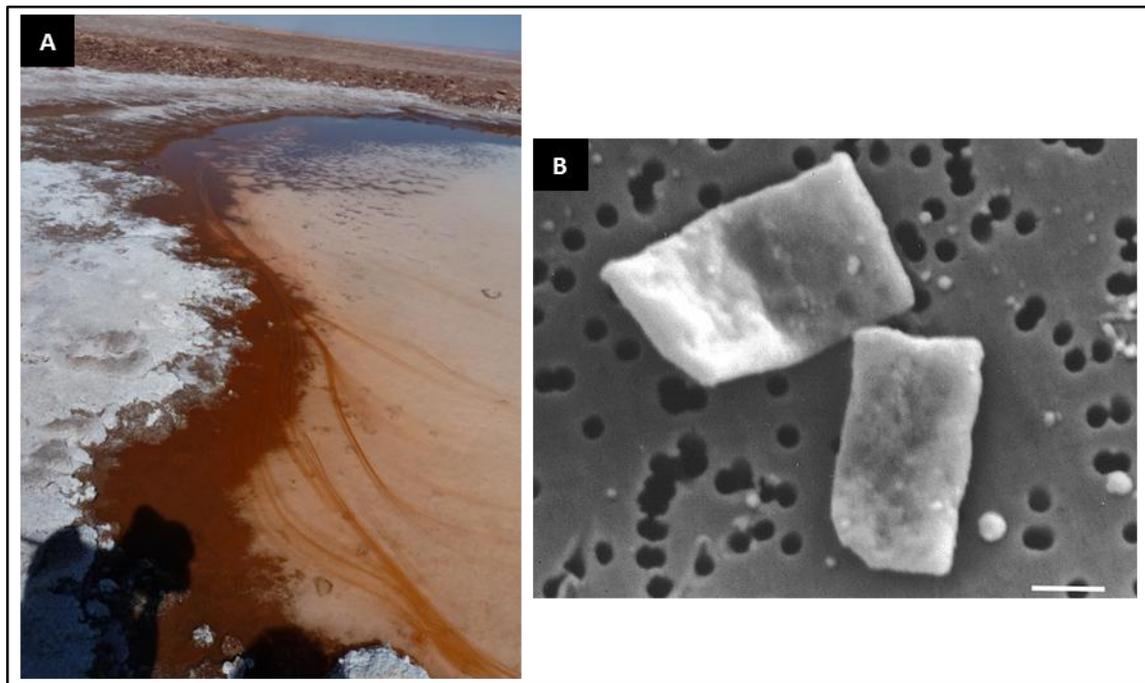


Figure 11. Archées halophiles. A : Salar du Lac d'Atacama au Chili. Les trainées rouges sont des colonies d'archées halophiles filamenteuses. Photographie de Purificación López-García. B : Haloquadratum walsbyi, une halobactériale à la forme caractéristique carrée et plate. Photographie de Francisco Rodríguez-Valera. Barre= 1µm

Acidophiles

Enfin, des archées sont adaptées aux milieux extrêmement acides, particulièrement les *Thermoplasmatales*, *Aciduliprofundum boonei*, archée du groupe DHVE2, phylogénétiquement proches des *Thermoplasmatales*, et les archées du groupe ARMAN. La plupart des *Thermoplasmatales* et les ARMAN ont été mises en évidence dans des mines acides ou dans des champs hydrothermaux volcaniques (Figure 12), à des pH compris entre 1 et 2 (Dopson et al. 2004; Darland et al. 1970). L'organisme *Picrophilus torridus*, isolé dans un sol volcanique sulfureux au Japon, est capable de croître à un pH de 0 (C Schleper et al. 1995). L'ordre *Thermoplasmatales* est aussi caractérisé par l'absence de paroi cellulaire (Darland et al. 1970; Dopson et al. 2004; C Schleper et al. 1995), leur donnant ainsi une morphologie proche des bactéries du groupe *Mycoplasmatales* avec lesquels ces organismes étaient classés jusqu'à ce que Woese analyse leurs ARNr SSU (C R Woese, Magrum, and Fox 1978). Ces organismes acidophiles peuvent aussi être thermophiles, vivant à des températures comprises entre 35 et 70°C. *Aciduliprofundum boonei* a été isolé depuis des échantillons de cheminées hydrothermales sous-marines. Cet organisme vit à un pH légèrement plus élevé (entre 3 et 4) mais à des températures plus hautes, autour de 70°C (Reysenbach et al. 2006).



Figure 12. Mine acide à Rio Tinto en Espagne. La couleur rouge-brune de l'eau sur les photos de gauche est due à la présence de métaux, et particulièrement de fer oxydé en solution. L'eau de la photo de droite est à pH très bas, autour de 1, et sa couleur verte est due à du fer ferreux, non oxydé. Photographie de Purificación López-García.

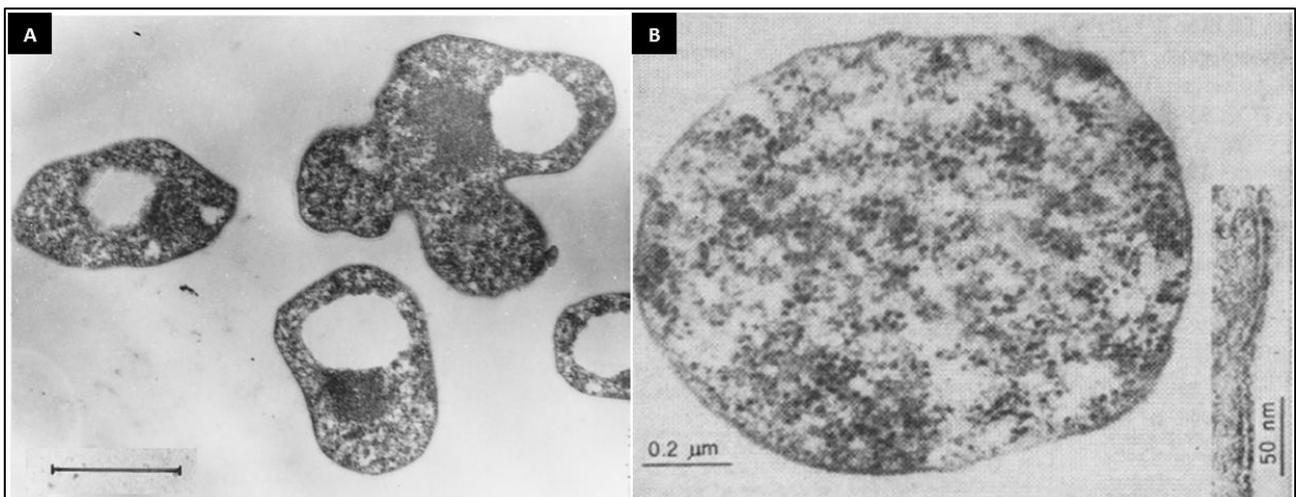


Figure 13. Archées acidophiles. A : *Picrophilus oshimae* en coupe fine. Cette archée a été isolée dans une source chaude sulfureuse au Japon, à une température de l'ordre de 55°C et un pH inférieur à 0,5 (C Schleper et al. 1995) Barre=1μm. B : *Thermoplasma acidophilum* en coupe fine à gauche. Elargissement sur la membrane à droite, montrant l'absence de paroi cellulaire, caractéristique des Thermoplasmatales (Darland et al. 1970).

Méthanogènes et méthanotrophes

Les méthanogènes sont des organismes capables de produire du méthane à partir de composés organiques ou inorganiques (CO₂, méthanol ou composés méthylés), le méthane étant le produit principal de ce métabolisme énergétique. Seules des archées possèdent les enzymes nécessaires à ce métabolisme complexe (Offre, Spang, and Schleper 2013; Y. Liu and Whitman 2008). Actuellement sept ordres de méthanogènes ont été décrits, tous des euryarchées, bien que ces ordres ne forment pas un ensemble monophylétique. Les *Methanopyrales*, les *Methanobacteriales* et les *Methanococcales* sont plus proches entre eux et ont été proposés comme « Methanogen Class I » par Bapteste et collaborateurs en 2005 (Eric Bapteste, Brochier, and Boucher 2005). Les *Methanosarcinales*, les *Methanomicrobiales* et les *Methanocellales* (anciennement Rice Cluster I (Sakai et al. 2008)) sont plus proches entre eux et ont été proposés comme « Methanogen Class II » (Eric Bapteste, Brochier, and Boucher 2005). Cependant, la monophylie de ces deux « classes » est toujours discutée (Brochier-Armanet, Forterre, and Gribaldo 2011). Le septième ordre de méthanogènes a été proposé très récemment par Paul et collaborateurs (Paul et al. 2012) suite à l'analyse de nombreuses séquences environnementales se plaçant à proximité des *Thermoplasmatales* (Mihajlovski, Alric, and Brugère 2008; Paul et al. 2012). Le gène *mrcA* qui code pour la sous-unité alpha du méthyl-coenzyme M, marqueur moléculaire classique des archées méthanogènes, a été amplifié à partir des mêmes échantillons que les ARNr SSU étudiés. Les phylogénies de ces deux marqueurs (*mrcA* et ARNr SSU) correspondent, renforçant l'idée que ces séquences environnementales représenteraient en réalité un nouvel ordre d'archées méthanogènes, et le nom de *Methanoplasmatales* a été proposé. Ces organismes ont été observés dans des milieux très divers, mais toujours anoxiques, tels que des sédiments marins, d'eau douce, des sols de marais, des rizières, le rumen de bovins, les tractus gastro-intestinaux de différents animaux (gros mammifères, insectes...), des tourbières, mais aussi dans des milieux considérés comme plus extrêmes, comme des sources chaudes, des hydrocarbures pétroliers ou des sédiments hypersalins (Y. Liu and Whitman 2008).

Les groupes ANME-1, -2 et -3 sont des groupes d'euryarchées non cultivées méthanotrophes, détectées dans les mêmes environnements que les méthanogènes (Knittel and Boetius 2009; Mills et al. 2005; Offre, Spang, and Schleper 2013). Ces organismes sont caractérisés par leur capacité à oxyder le méthane en anaérobiose. Il est donc logique qu'ils vivent dans les mêmes milieux que les méthanogènes. La plupart du temps des bactéries capables de réduire le sulfate se trouvent aussi dans ces mêmes milieux, ce qui permettrait une oxydation anaérobique du méthane (« AOM » pour « Anaerobic Oxydation of Methane ») dépendante du sulfate avec production d'énergie pour les archées ANME et les bactéries sulfato-réductrices (Offre, Spang, and

Schleper 2013; Knittel and Boetius 2009). D'un point de vue environnemental, ce processus AOM balancerait la production de méthane biologique et son émission dans l'atmosphère puisqu'il consommerait une grande part de cette production, particulièrement dans les environnements marins (Reeburgh 2007). L'impact écologique de ces deux populations d'archées, méthanogènes et méthanotrophes, est donc très important.

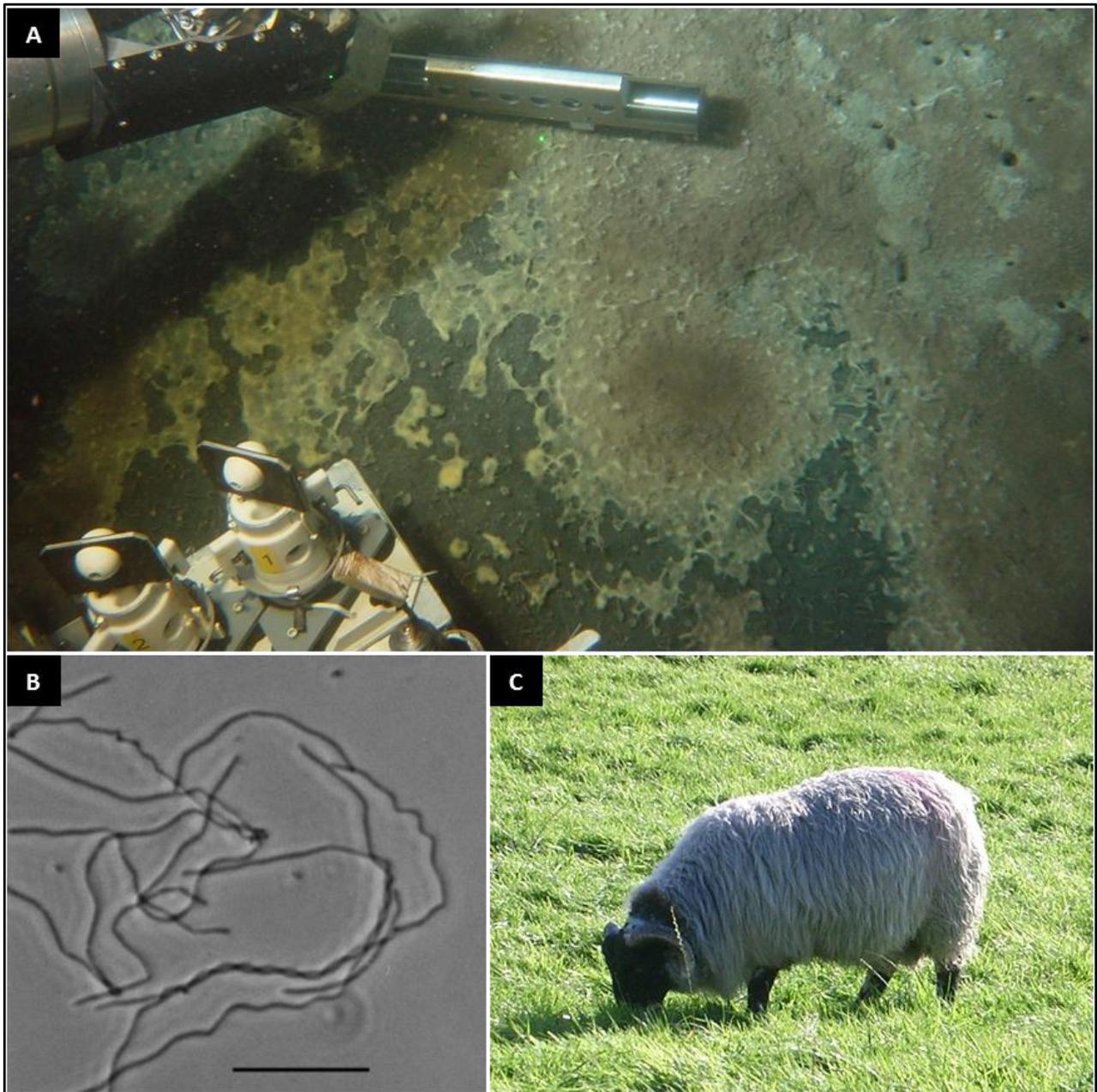


Figure 14. Archées méthanogènes et leurs habitats. A : Tapis microbien anoxique sous-marin, dans lequel sont présentes des archées méthanogènes et des méthanotrophes. Photographie de Purificación López-García. B : *Methanobacterium oryzae*, une archée méthanogène isolée dans un champ de riz aux Philippines. Barre= 10µm. (Jouliau et al. 2000). C : *Ovis arie*. Certains méthanogènes ont été isolés depuis leur tractus digestif, comme *Methanosarcina barkeri*. Photographie de Celine Petitjean.

b. Les Nanoarchaeota

Actuellement *Nanoarchaeum equitans* est le seul représentant cultivé de ce phylum et dont le génome ait été séquencé (H. Huber et al. 2002). C'est une archée hyperthermophile, de toute petite taille (~ 400 nm), avec un génome très réduit et vivant justement en symbiose avec *Ignicoccus hospitalis*, une crenarchée avec laquelle elle a pu être cultivée (H. Huber et al. 2002) (Figure 13). Le séquençage de son génome a révélé des caractéristiques particulières typiques des organismes symbiotiques (Waters et al. 2003), à savoir un génome très réduit avec absence de beaucoup de gènes essentiels du métabolisme (biosynthèse des acides aminés, nucléotides, cofacteurs et lipides), des gènes intervenant dans les différentes voies d'assimilation du carbone ou de trois ARNt (glutamate, histidine et tryptophane). L'absence de ces gènes rend *N. equitans* complètement dépendant d'*I. hospitalis* et tend à confirmer son évolution par réduction de génome due à un mode de vie symbiotique. On trouve aussi 10 gènes coupés en deux dans ce génome alors qu'ils sont codés en une seule séquence chez la majorité des autres archées. Ce nombre est assez important comparativement aux autres archées, et certains des gènes concernés sont impliqués dans des processus cellulaires essentiels, comme l'ADN polymérase I, l'alanyl-ARNt synthétase ou la sous-unité β de l'ARN polymérase. Sa position phylogénétique et la justification de la création d'un phylum à part entière seront discutées plus tard.

Très récemment un deuxième génome de Nanoarchaeota a été séquencé par métagénomique : « Nst1 » (Podar et al. 2013). Cet organisme vivrait en symbiose avec une crenarchée de l'ordre des *sulfolobales*. Bien que son génome soit environ deux fois plus grand (de l'ordre de 1 mégabase) que celui de *N. equitans* (de l'ordre de 0,5 mégabase), on retrouve aussi des pertes massives de gènes du métabolisme et un certain nombre de gènes coupés en deux, certains communs avec *N. equitans*. Ce génome semble donc avoir subi une évolution par réduction mais moins drastique que celui de *N. equitans*.

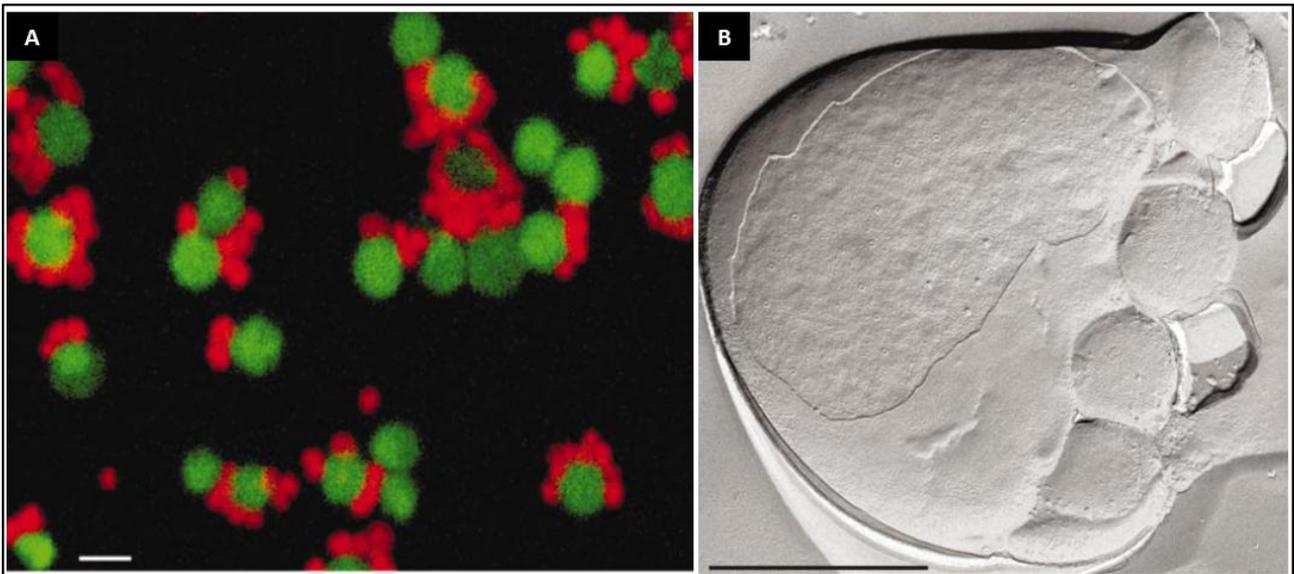


Figure 15. Nanoarchaeum equitans et Ignicoccus hospitalis, son hôte. (H. Huber et al. 2002). Barre= 1 μ m. A : Image de microscopie confocale laser après hybridation avec le marqueur rouge CY3 pour Nanoarchaeum equitans et le marqueur vert rhodamine-green pour I. hospitalis. B : Image de microscopie électronique d'I. hospitalis et de quatre cellules de N. equitans accrochées.

c. Les Crenarchaeota

Les Crenarchaeota sont des organismes hyperthermophiles ou thermophiles répartis en cinq ordres : les *Desulfurococcales*, les *Sulfolobales*, les *Thermoproteales*, ainsi que les *Acidilobales* et les *Fervidicoccales*, deux nouveaux ordres proposés en 2009 (Prokofeva et al. 2009) et en 2010 (Perevalova et al. 2010). La plupart des *Thermoproteales* et des *Desulfurococcales* sont anaérobies et hyperthermophiles, alors que les *Sulfolobales* ont plutôt tendance à être aérobies et que leurs températures de croissance sont plus basses que celles observées dans les deux autres groupes. Les environnements dans lesquels on les trouve peuvent être les mêmes que ceux dans lesquels sont présents les euryarchées hyperthermophiles (Figure 8), à savoir des sources chaudes terrestres ou des sources hydrothermales sous-marines, de surface ou profondes, par exemple les fumeurs noirs et blancs présents au niveau des dorsales océaniques (Figure 16).

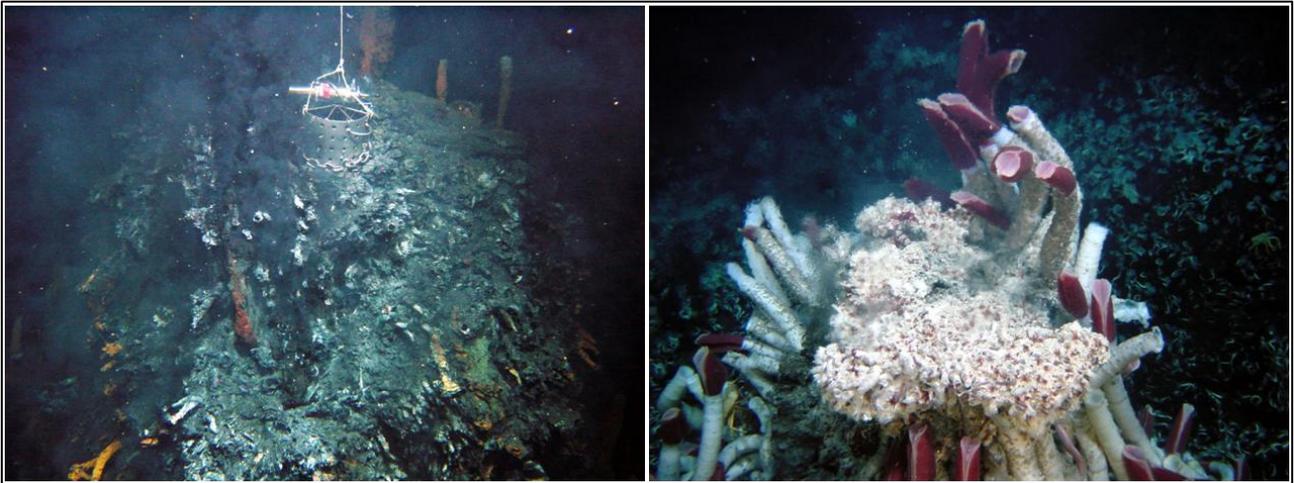


Figure 16. Fumeur noir de la dorsale océanique East Pacific Rise. Environnement hyperthermophile dans lequel se trouvent particulièrement des *Desulfurococcales*. Photographie de Purificación López-García.

Les *Desulfurococcales* sont les crenarchées les plus hyperthermophiles. Leurs températures de croissance varient autour de 90°C, et peuvent monter jusqu'à plus de 100°C. Par exemple, *Pyrolobus fumarii* a une température optimale de croissance de 106°C mais est capable de croître jusqu'à 113°C (Anderson, Göker, et al. 2011). Cette archée a été isolée dans des fumeurs hydrothermaux des fonds sous-marins, l'augmentation de la pression avec la profondeur permettant de telles températures. La majorité des *Desulfurococcales* sont anaérobies, stricts ou facultatifs, et leurs métabolismes sont assez variés. Par exemple, *Hyperthermus butylicus* (Brügger et al. 2007) ou *Ignicoccus hospitalis* (Podar et al. 2008) sont anaérobies et sulfato-réductrices et ont été isolées dans des sources hydrothermales sulfureuses. Certaines sont hétérotrophes, comme *Staphylothermus hellenicus* (Anderson, Wirth, et al. 2011), *Thermosphaera aggregans* (Spring et al. 2010), ou *Ignisphaera aggregans* (Göker et al. 2010). *Desulfurococcus fermentans*, isolée dans une source chaude au Kamchatka en Russie, est même capable de dégrader la cellulose (Susanti et al. 2012). *Ignicoccus hospitalis*, l'hôte de *Nanoarchaeum equitans*, est une désulfurococcale isolée à partir d'une source chaude sous-marine en Islande. Elle a une température de croissance de 95°C (Podar et al. 2008). Parmi les membres du genre *Ignicoccus*, seul *I. hospitalis* est l'hôte de *N. equitans* ou d'un symbiote tout court. Les *Ignicoccus* ont une cytologie particulière avec une double membrane cellulaire et un périplasme très important entre la membrane cellulaire interne et la membrane externe, dans lequel de nombreuses vésicules sont présentes (Rachel et al. 2002).

Les *Acidilobales* ont été proposés récemment par Prokofeva et collaborateurs en 2009 (Prokofeva et al. 2009), et incluraient les organismes des genres *Acidilobus* et *Caldisphaera*. Ces organismes sont thermophiles, anaérobies, organotrophes et acidophiles. Ils ont été isolés dans différentes sources chaudes acides, au Kamchatka en Russie, à Laguna aux Philippines ou dans le parc de Yellowstone aux Etats-Unis. Les analyses phylogénétiques de leur ARNr SSU les placent

groupe frère des *Desulfurococcales*, conduisant les auteurs à proposer un nouvel ordre de Crenarchaeota. Le génome complet d'*Acidilobus saccharovorans*, isolée au Kamchatka, a été séquencé en 2009 (Prokofeva et al. 2009). Il a révélé de nombreuses voies métaboliques laissant supposer qu'il est capable d'utiliser une très large gamme de composés organiques produits par les organismes lithoautotrophes vivant dans les mêmes environnements.

Peu de temps après, en 2010, un nouvel ordre de Crenarchaeota est proposé, les *Fervidicoccales* (Perevalova et al. 2010), après l'isolement de *Fervidicoccus fontis* dans les sources chaudes volcaniques du Kamchatka. C'est un organisme thermoacidophile (croissance optimale entre 65 et 70°C, à pH 5,5-6), anaérobie et organotrophe. L'analyse phylogénétique de l'ARNr SSU de cet organisme le place en groupe frère des *Acidilobales*, avec d'autres séquences environnementales trouvées dans des environnements de type sources hydrothermales à Yellowstone ou en Islande. Cette position phylogénétique a conduit les auteurs à proposer que ces séquences soient les représentantes d'un nouvel ordre.

Les *Thermoproteales* sont aussi des organismes hyperthermophiles et anaérobies pour la plupart. Ils vivent aussi dans ces mêmes milieux et principalement dans des sources chaudes terrestres sulfureuses, les solfatares, à des températures légèrement plus basses que les *Desulfurococcales*. Certains sont légèrement acidophiles et ont été isolés à des pH inférieurs à 7, comme *Thermoproteus tenax*, isolé à un pH de 5,6 dans un solfatare en Islande (Zillig et al. 1981; Siebers et al. 2011). On trouve aussi une variété importante de métabolismes : des organismes organotrophes comme *Caldivirga maquilingensis*, hétérotrophe et microaérophile, des organismes capables de croître de façon lithotrophe et organotrophe, comme *Pyrobaculum arsenaticum*, utilisant l'hydrogène comme donneur d'électrons pour la réduction de l'arsenate et de composés sulfurés pour une croissance chimiolithoautotrophe et qui utilise ces derniers composés comme accepteurs d'électrons pour une croissance organotrophe (R. Huber et al. 2000).

Enfin les *Sulfolobales* sont des organismes thermophiles ou hyperthermophiles plus modérés que les autres ordres, leurs températures de croissance sont plutôt comprises entre 60 et 85°C, et on les trouve principalement dans des sources chaudes terrestres, comme celles de Solfatare en Italie, de Yellowstone ou du Kamchatka. On trouve aussi plus d'organismes aérobies dans ce groupe, alors qu'ils sont minoritaires dans les autres ordres de Crenarchaeota, et d'acidophiles extrêmes, vivant à des pH entre 2 et 5. Par exemple *Sulfolobus acidocaldarius*, le premier hyperthermoacidophile caractérisé depuis un solfatare terrestre par Brock et collaborateurs en 1972 (Brock et al. 1972), croît entre 75 et 80°C, à un pH entre 2 et 3, et est strictement aérobie et organotrophe (Chen et al. 2005).

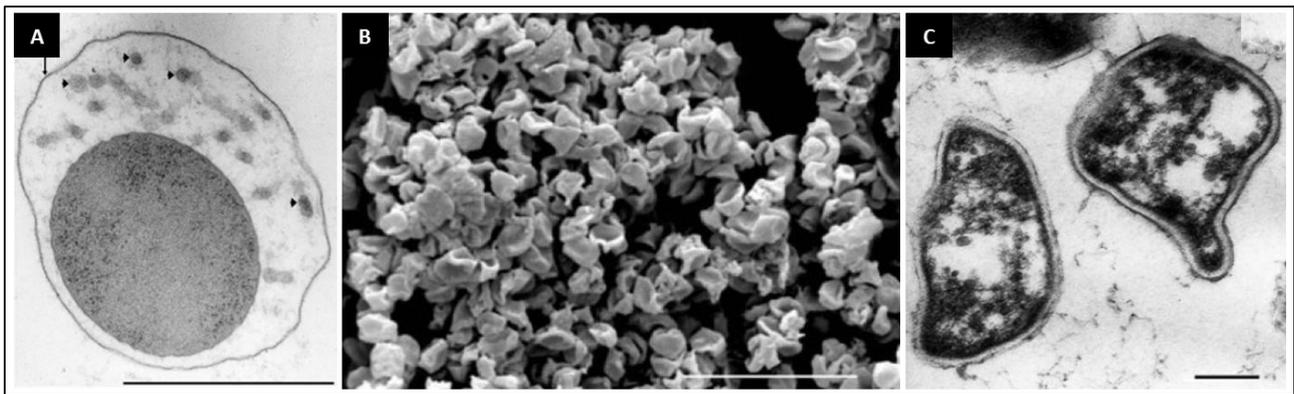


Figure 17. Crenarchées. A : Image de microscopie électronique d'une cellule d'Ignicoccus (Rachel et al. 2002). On observe le large périplasm entre la membrane interne et la membrane externe. Barre=1 μ m. B : Colonie de cellules de Metallosphaera cuprina, une sulfobactérie. Barre= 5 μ m. (L.-J. Liu et al. 2011). C : Coupe fine de cellules d'Acidilobus saccharovorans, une acidilobale. Barre= 0,5 μ m. (Prokofeva et al. 2009).

d. Les Thaumarchaeota

Les Thaumarchaeota ont été découverts par DeLong et Fuhrmann en 1992 et nommés « groupe I ». Phylogénétiquement proches des Crenarchaeota elles étaient considérées comme des « Crenarchaeota mésophiles » (DeLong 1992; Fuhrman, McCallum, and Davis 1992). En 2008, après la publication du génome de *Cenarchaeum symbiosum*, Brochier-Armanet et collaborateurs ont réanalysé la position de ce groupe par des méthodes de phylogénomique et ont proposé que les membres du groupe I ne soient pas des Crenarchaeota mais un phylum à part entière. Ils ont proposé le nom de « *Thaumarchaeota* » (Brochier-Armanet et al. 2008). De nombreux organismes se placent parmi les Thaumarchaeota, mais très peu sont cultivés axéniquement, le premier a été *Nitrosopumilus maritimus* (Figure 18), 18 ans après la découverte de ces archées (Könneke et al. 2005).

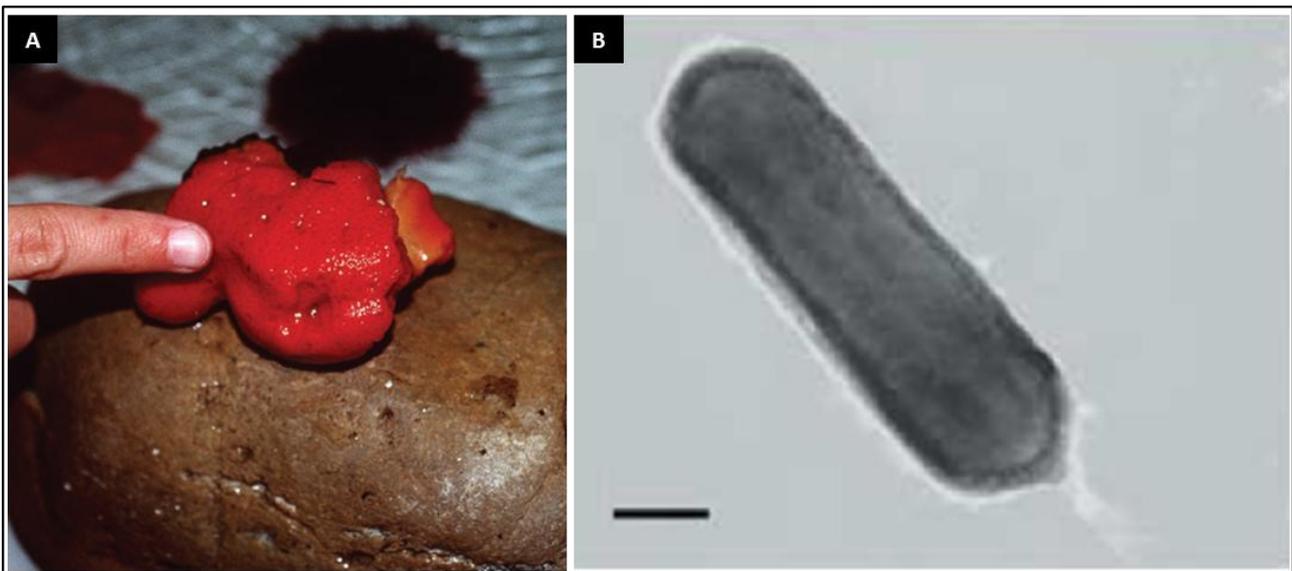


Figure 18. Thaumarchées. A : *Axinella mexicana*, une éponge, est l'hôte de *Cenarchaeum symbiosum* (DeLong 2003). B : *Nitrosopumilus maritimus*, première thaumarchée isolée en 2005, vivant dans les milieux marins (Könneke et al. 2005) Barre=0,1µm.

Les Thaumarchaeota sont des organismes capables de réaliser la première étape de la nitrification, à savoir l'oxydation aérobie de l'ammoniac, NH_3 , en nitrite, NO_2^- . Cette capacité est très rare dans le vivant, et jusqu'à la description des thaumarchées elle n'était connue que chez quelques beta- et gamma-protéobactéries (AOB, pour Ammonia-Oxydizing Bacteria). Cette réaction est catalysée par l'ammonia monooxygénase AMO, composée de trois sous-unités : AmoA, AmoB et AmoC chez les archées (Stahl and de la Torre 2012). Les thaumarchées dominent une partie de la biosphère des océans en eau profonde, comme l'ont montré Karner et collaborateur en 2001 (Karner, DeLong, and Karl 2001). Leur capacité d'oxydation de l'ammoniac en nitrite associée à l'importance en biomasse de ces organismes en fait des acteurs majeurs du cycle biogéochimique de l'azote.

Depuis, de nombreuses études environnementales ont démontré leur présence dans la plupart des environnements. En effet, bien qu'elles aient été découvertes dans des milieux mésophiles, des thaumarchées sont présentes dans des milieux extrêmement divers, et pas uniquement mésophiles. Dans l'océan profond, bien sûr, ou moins profond, par exemple *Cenarchaeum symbiosum* qui vit en symbiose avec l'éponge *Axinella mexicana* (Preston et al. 1996). On les trouve aussi dans les sols (*Nitrososphaera viennensis* (Tourna et al. 2011)), dans des sédiments marins ('*Candidatus Nitrosopumilus salaria*' (Mosier et al. 2012) ou '*Candidatus Nitrosopumilus sediminis*' (Park et al. 2012)). D'autres lignées sont présentes dans des milieux thermophiles tels que des tapis microbiens de sources chaudes (*Nitrososphaera gargensis* (Hatzenpichler et al. 2008)) ou dans des sources chaudes à Yellowstone ('*Candidatus Nitrosocaldus yellowstonii*' (de la Torre et al. 2008)) ; ces deux

derniers organismes sont donc thermophiles. Récemment de nouveaux groupes très divergents, HTC1 et HTC2, pour « Hot Thaumarchaeota-related Clade » 1 et 2, ont été décrits par Eme et collaborateurs en 2013 (Eme et al. 2013) à partir d'échantillons de sources chaudes du Kamchatka. La [Figure 19](#) montre une phylogénie de l'ARNr SSU des Thaumarchaeota issue de cet article, représentative de leur diversité.

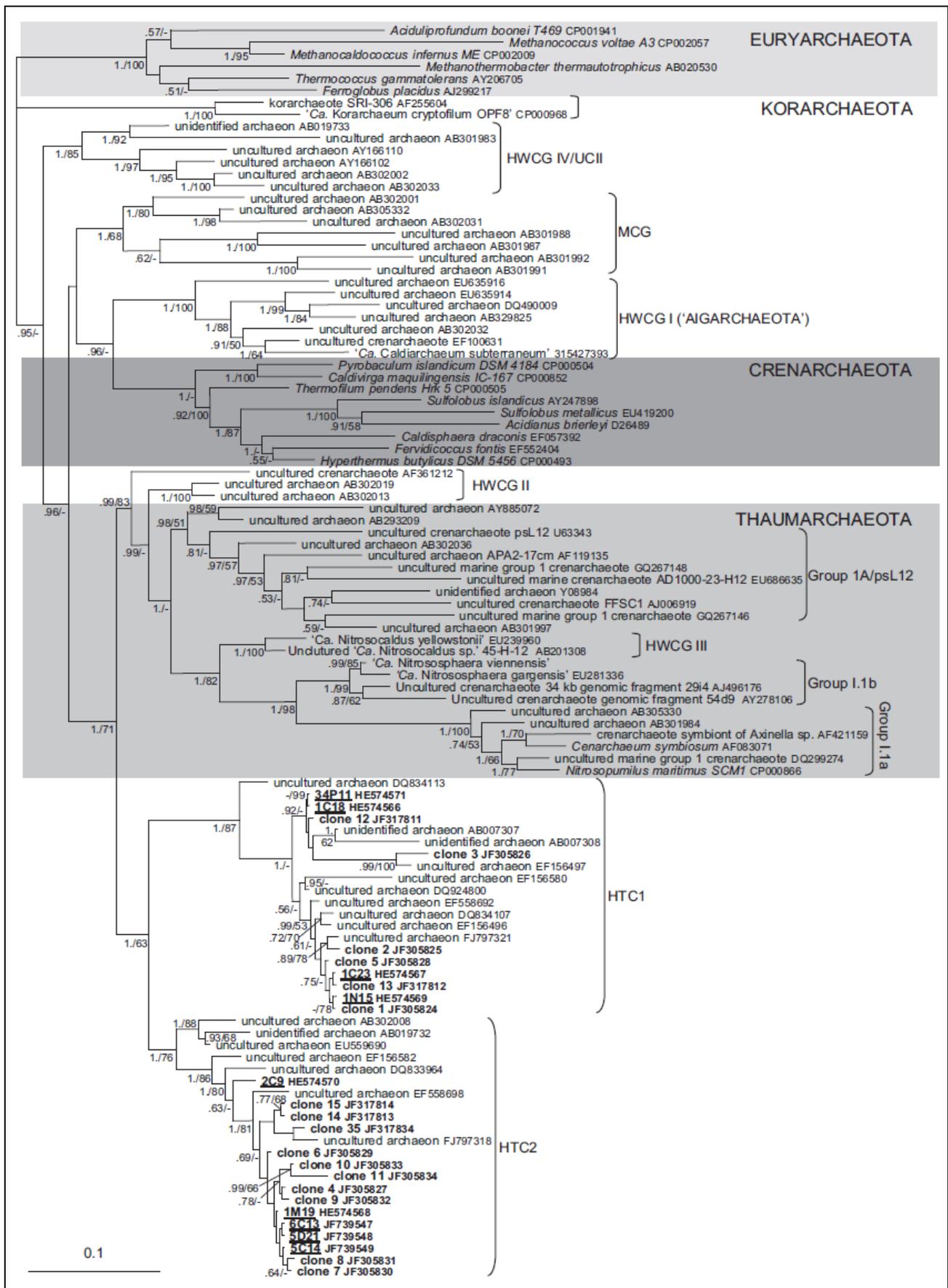


Figure 19. Phylogénie des Archées centrée sur les Thaumarchaeota non cultivées. Arbre non raciné, inféré sur 1064 positions de 105 séquences d'ARNr SSU (Eme et al. 2013).

Il est intéressant de noter que des organismes mésophiles et thermophiles sont présents chez les Thaumarchaeota mais que les organismes mésophiles sont polyphylétiques, plusieurs transitions ont donc dû se produire de la thermophilie à la mésophilie ou inversement. Plusieurs travaux de López-García et collaborateurs ont montré que les transferts horizontaux de gènes ont joué un rôle important dans l'évolution du génome des Thaumarchaeota marines (groupe I) et particulièrement dans l'adaptation à leur milieu océanique profond froid et pauvre en nutriments (López-garcía, Brochier, Moreira, & Rodríguez-valera, 2004 ; Brochier-Armanet et al., 2011) (cf. Chapitre 3).

La Figure 19 illustre également le fait que la majeure partie de la diversité des thaumarchées est non cultivée. Ainsi, très peu de génomes complets sont disponibles. En dehors de celui de *Cenarchaeum symbiosum* en 2006 (Hallam et al. 2006), de plus en plus de génomes sont séquencés pour ce phylum ces deux dernières années, mais la plupart appartiennent au groupe I.1a, proche de *Nitrosopumilus maritimus*, séquencé en 2010 (Walker et al. 2010).

e. Les Aigarchaeota

Ce phylum proposé par Nunoura et collaborateurs en 2011 (Nunoura et al. 2011) correspond au groupe HWCG I, visible sur la Figure 19, où il se place en groupe frère des Crenarchaeota. Ce phylum a été proposé en même temps que la publication du génome de '*Candidatus Caldiarchaeum subterraneum*' dont il a déjà été question plus tôt. Ces organismes sont présents dans des environnements thermophiles classiques : sources chaudes terrestres (Nunoura et al. 2005) et sous-marines, sources hydrothermales de mer profonde (Reykjaví et al. 2001). Aucun représentant n'est cultivé. Le génome de '*Ca. Caldiarchaeum subterraneum*' a été obtenu à partir de l'enrichissement d'un tapis microbien de surface dans de l'eau géothermale d'une mine d'or, dans lequel ce groupe HWCGI était dominant, et a été reconstruit à partir d'une banque de données métagénomiques (Nunoura et al. 2011). Dans cette publication, les analyses phylogénétiques placent '*Ca. Caldiarchaeum subterraneum*' en groupe frère des Thaumarchaeota.

f. Les Korarchaeota

Le phylum Korarchaeota a été proposé en 1996 par Barns et collaborateurs (Barns et al. 1996) sur la base de séquences environnementales produites à partir d'une source chaude de Yellowstone. En 2006 Auchtung et collaborateurs ont étudié la diversité des Korarchaeota en construisant des amorces ciblant spécifiquement l'ARNr SSU de ce groupe. Ils ont échantillonné un grand nombre d'environnements hydrothermaux, des sources chaudes à Yellowstone et des fumeurs

hydrothermaux de fonds sous-marins, mais aussi d'environnements non hydrothermaux tels que du compost, des végétaux, de la salive humaine, des sols, de l'eau, des sédiments de rivières, etc. Le résultat de cette étude est clair : les korarchées sont détectées uniquement dans des environnements thermophiles (température supérieure à 55°C). La plupart des nouvelles séquences, huit, ont été détectées dans les échantillons provenant des différentes sources chaudes de Yellowstone, et une autre séquence dans un échantillon de cheminée hydrothermale sous-marine. Les séquences provenant d'études précédentes sont aussi issues d'environnements thermophiles et l'ensemble forme un groupe monophylétique robuste, branchant avant la divergence entre Euryarchaeota et Crenarchaeota (Auchtung, Takacs-Vesbach, and Cavanaugh 2006). Cette étude confirme la présence de représentants de ce groupe phylogénétique dans divers environnements, mais toujours à haute température. Le seul génome de Korarchaeota disponible à ce jour a été publié en 2008. Il s'agit du génome de '*Ca. Korarchaeum cryptofilum*', reconstruit à partir d'enrichissement de cellules provenant de la source chaude « Obsidian Pool » à Yellowstone (Elkins et al. 2008). Ce génome se révèle assez particulier, avec l'absence de nombreux gènes permettant la synthèse *de novo* de molécules importantes, telles que les purines, le coenzyme A ou d'autres cofacteurs, présents chez la plupart des autres archées. Il partage avec les Euryarchaeota des fonctions liées à la maturation des ARNt, à la réplication et réparation de l'ADN et à la division cellulaire, mais la composition protéique de son ribosome et de son ARN Polymérase est plus proche de celles des Crenarchaeota. Sa position phylogénétique et ses relations avec ces deux groupes peuvent donc être considérées comme cohérentes et ces caractéristiques pourraient avoir été présentes dans leur ancêtre commun (Elkins et al. 2008). Ce groupe reste aujourd'hui encore très peu connu et il serait intéressant d'avoir plus de données, génomiques et de culture, sur ces organismes ayant une place si importante dans l'histoire évolutive des archées.

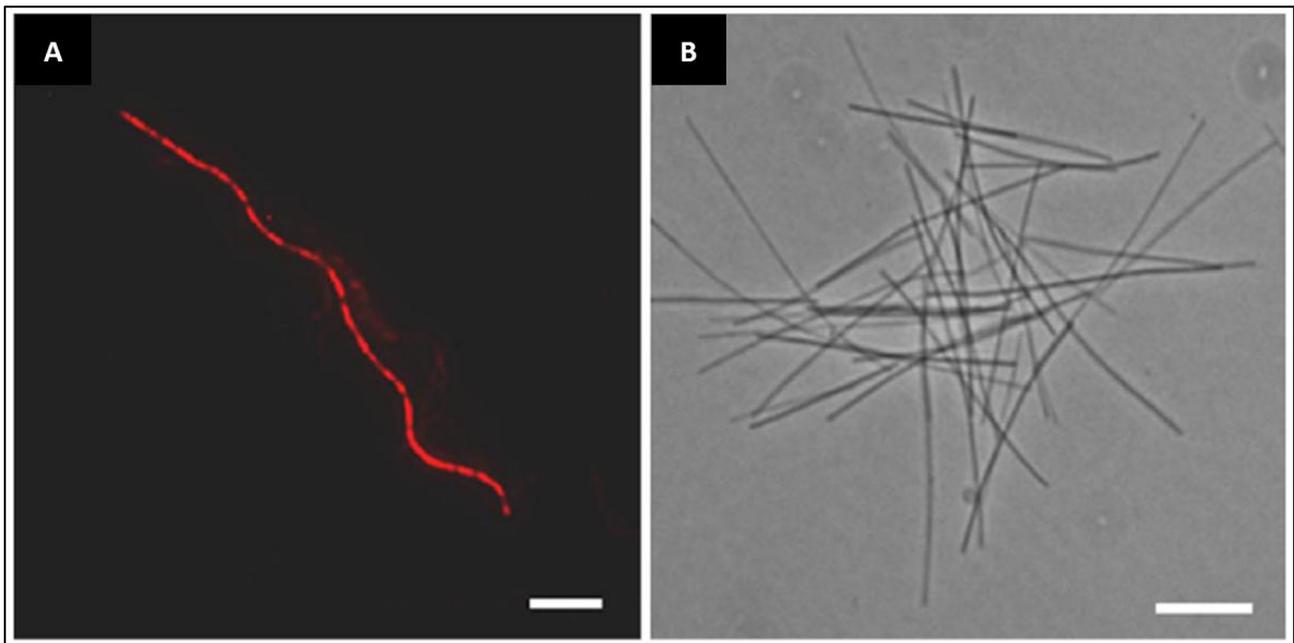


Figure 20. '*Candidatus Korarchaeum cryptofilum*'. Ces deux photographies de microscopie montrent bien la forme filamenteuse des cellules de '*Ca. Korarchaeum cryptofilum*'. Barre= 5 μ m. A : Analyse en FISH. B : Enrichissement de '*Ca. Korarchaeum cryptofilum*'. (Elkins et al. 2008).

g. Conclusion

La [Figure 7](#), une phylogénie construite à partir de l'ARNr SSU, représente bien la diversité des archées, cultivées ou non, avec des génomes séquencés ou non. Les différents ordres d'Euryarchaeota et de Crenarchaeota apparaissent bien de façon monophylétique, ainsi que les différents groupes de Thaumarchaeota. Malgré tout, les relations entre ces groupes, ordres et phyla, ne sont pas résolues. Ceci montre que l'utilisation d'un seul marqueur, aussi puissant soit-il, n'est pas suffisante pour résoudre les relations anciennes entre autant d'organismes. Le séquençage massif de génomes a ouvert la porte à l'utilisation de nouveaux marqueurs, et surtout à l'utilisation conjointe de ces marqueurs dans des analyses phylogénomiques.

2. Phylogénies moléculaires inférées sur plusieurs marqueurs

a. La phylogénie des archées

Première phylogénie des archées

En 2002, Matte-Tailliez et collaborateurs analysent pour la première fois des concaténations de protéines ribosomiques pour inférer la phylogénie des archées (Matte-Tailliez et al. 2002). Précédemment, cette phylogénie n'avait été inférée qu'à partir de l'ARNr SSU. Comme je l'ai déjà évoqué plus tôt, ce marqueur est sujet à certains biais comme la différence de vitesse d'évolution entre séquences, ce qui conduit à de possibles attractions de longues branches (Philippe and Laurent 1998). Le nombre limité de sites qui le composent peut être aussi responsable d'une mauvaise résolution, particulièrement au niveau des relations anciennes entre organismes, à cause du phénomène de saturation mutationnelle (Simonetta Gribaldo and Brochier-Armanet 2009). De plus, la phylogénie de l'ARNr SSU reste la phylogénie d'un seul gène, et pas des organismes dans leur globalité. Bien que l'histoire évolutive de l'ARNr SSU puisse refléter celle des organismes, la confirmation de cette hypothèse par l'utilisation de nouveaux marqueurs est importante. Plusieurs protéines avaient été utilisées dans ce but, mais toujours dans des phylogénies individuelles avec un seul marqueur utilisé à chaque fois ; par exemple, la sous-unité B de l'ARN polymérase (Klenk and Zillig 1994), différentes aminoacyl-ARNt synthétases (C R Woese et al. 2000), ou la chaperonne HSP70/DnaK (R. S. Gupta 2000).

Les protéines ribosomiques sont extrêmement conservées au niveau fonctionnel et sont universelles (présentes dans l'ensemble du vivant), deux qualités qui avaient fait de l'ARNr SSU un marqueur de choix pour la phylogénie des organismes. Treize génomes d'archées étaient disponibles en 2002 : trois Crenarchaeota et dix Euryarchaeota. Les auteurs ont donc analysé la phylogénie de 53 protéines ribosomiques clairement orthologues dans ces génomes, les ARNr SSU et LSU et la sous-unité B de l'ARN polymérase (Matte-Tailliez et al. 2002). A partir de ces 53 protéines, les auteurs ont réalisé différentes concaténations ([Figure 21 A et B](#)) d'une part, et d'autre part, fait une analyse de chaque phylogénie individuelle afin de détecter de potentiels transferts horizontaux de gènes. Enfin, ils ont ciblé dans une analyse plus fine les gènes qui apparaissaient comme ayant pu être transférés. Ceci aura permis d'identifier huit protéines porteuses d'un signal divergent. Cette méthode avait été mise en place et utilisée pour l'étude de la phylogénie des bactéries à partir de marqueurs multiples la même année (Brochier et al. 2002).

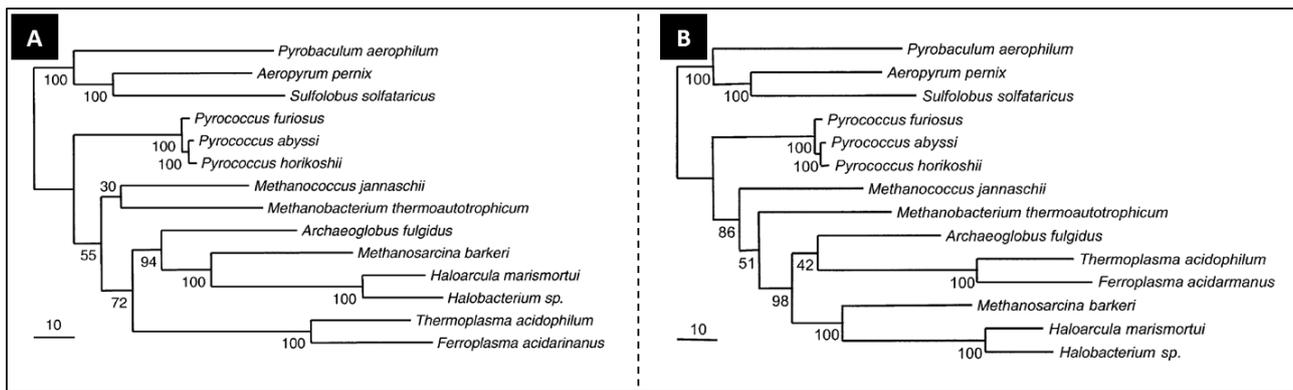


Figure 21. Phylogénies des archées publiées en 2002. Arbres inférés par maximum de vraisemblance. A : à partir de la concaténation de l'ARNr SSU et de l'ARNr LSU (3933 positions). B : à partir d'une concaténation de 45 protéines ribosomiques (6249 positions). (Matte-Tailliez et al. 2002).

La concaténation obtenue à partir de l'ensemble des 45 protéines ribosomiques ne montrant pas de transferts (Figure 21 B) a donné un résultat très similaire à celui obtenu par l'analyse de la concaténation des ARNr SSU et LSU (Figure 21 A) mais globalement avec un meilleur support des branches. Une dernière phylogénie a été inférée à partir de la sous-unité B de l'ARN polymérase. Elle a une topologie très proche de celles obtenues à partir des protéines ribosomiques. Le travail de Matte-Tailliez et collaborateurs montre que les protéines ribosomiques portent dans l'ensemble un signal phylogénétique très proche de celui de l'ARNr SSU et de la sous-unité B de l'ARN polymérase. Les différents marqueurs semblent donc partager la même histoire évolutive et les implications de cette découverte sont importantes. Tout d'abord, la congruence de l'histoire évolutive de ces différents systèmes permet donc de poser l'hypothèse que cette histoire soit représentative de l'histoire évolutive des organismes. D'autre part, les travaux de Woese sur la phylogénie des archées basée sur l'ARNr SSU sont donc globalement confirmés par l'analyse de ces nouveaux marqueurs, et l'ARNr SSU semble donc un bon marqueur pour la phylogénie des organismes, même s'il n'apporte pas une résolution optimale. Enfin il apparaît que si certaines protéines ribosomiques semblent avoir subi des transferts au cours de leur histoire évolutive, l'analyse de leurs phylogénies individuelles comparées à une phylogénie de référence inférée à partir de plusieurs marqueurs permet de les écarter.

Les deux arbres montrés dans la Figure 21 A et B ont été racinés arbitrairement entre les Euryarchaeota et les Crenarchaeota, en accord avec les phylogénies universelles publiées précédemment. Seulement trois espèces de Crenarchaeota sont représentées pour dix d'Euryarchaeota. Parmi ces dernières, on peut observer 1) le groupe monophylétique formé par les trois *Pyrococcus*, première divergence au sein des Euryarchaeota ; 2) celui formé par deux *Halobacteriales* (*Haloarcula marismortui* et *Halobacterium sp.*) et *Methanosarcina barkeri*, seule représentante des *Methanosarcinales* ; et 3) celui formé par les *Thermoplasmatales* (*Thermoplasma*

acidophilum et *Ferroplasma acidarmanus*). Ces cinq espèces forment d'ailleurs un groupe monophylétique avec *Archaeoglobus fulgidus*, représentante des *Archaeoglobales*. Par contre, *Methanococcus jannaschii* (une méthanococcale) et *Methanobacterium thermoautotrophicum* (une méthanobactériale) sont groupe frère dans l'arbre inféré sur les ARNr (Figure 21 A) et paraphylétiques dans l'arbre inféré sur les protéines ribosomiques (Figure 21 B). Ces deux espèces branchent néanmoins toujours juste après les *Thermococcales*.

Evolution de la méthanogenèse et implications pour la phylogénie des archées

Peu de temps après, le génome de *Methanopyrus kandleri* est séquencé (Slesarev et al. 2002). Alors que les premières analyses phylogénétiques de son ARNr SSU en 1991 (Burggraf et al. 1991) le plaçaient à la base de l'arbre des archées, les analyses phylogénétiques présentées dans l'article de Slesarev et collaborateurs, basées sur des concaténations de protéines ribosomiques ou sur le contenu en gènes, placent cet organisme au sein des Euryarchaeota, après la divergence des *Thermococcales* et plus précisément, groupe frère des *Methanococcales* et des *Methanobacterales*. La Figure 22 illustre ces deux positions, bien qu'elle ne soit pas issue des articles en question. La position de *M. kandleri* au sein des archées est très importante pour comprendre l'évolution de la méthanogenèse. L'ensemble des méthanogènes connus partagent les mêmes enzymes homologues capables de réaliser la méthanogenèse (Eric Bapteste, Brochier, and Boucher 2005). Si *M. kandleri* a divergé avant les Euryarchaeota et les Crenarchaeota, il est possible que l'ancêtre des archées ait été méthanogène, et que cette voie métabolique ait été perdue dans les lignées non méthanogènes par la suite. L'hypothèse alternative est que les gènes impliqués dans la méthanogenèse auraient été transférés entre méthanogènes.

Deux articles, publiés en 2004 et 2005 par Brochier et collaborateurs (Brochier, Forterre, and Gribaldo 2004) et Bapteste et collaborateurs (Eric Bapteste, Brochier, and Boucher 2005) traitent de la position de *M. kandleri* chez les archées et de l'évolution de la méthanogenèse. Le premier (Brochier, Forterre, and Gribaldo 2004) est une analyse de la phylogénie des archées à partir d'une mise à jour des jeux de données de 53 protéines ribosomiques construits précédemment par Matte-Taille et collaborateurs en 2002 et d'un nouveau jeu de données construit à partir de 15 protéines impliquées dans la transcription (12 sous-unités de l'ARN polymérase ; A', A'', B, D, E', E'', F, H, K, L, N, P ; et les facteurs de transcription NusA, NuG et TFS). La Figure 22 montre les arbres obtenus à partir de ces deux jeux de données.

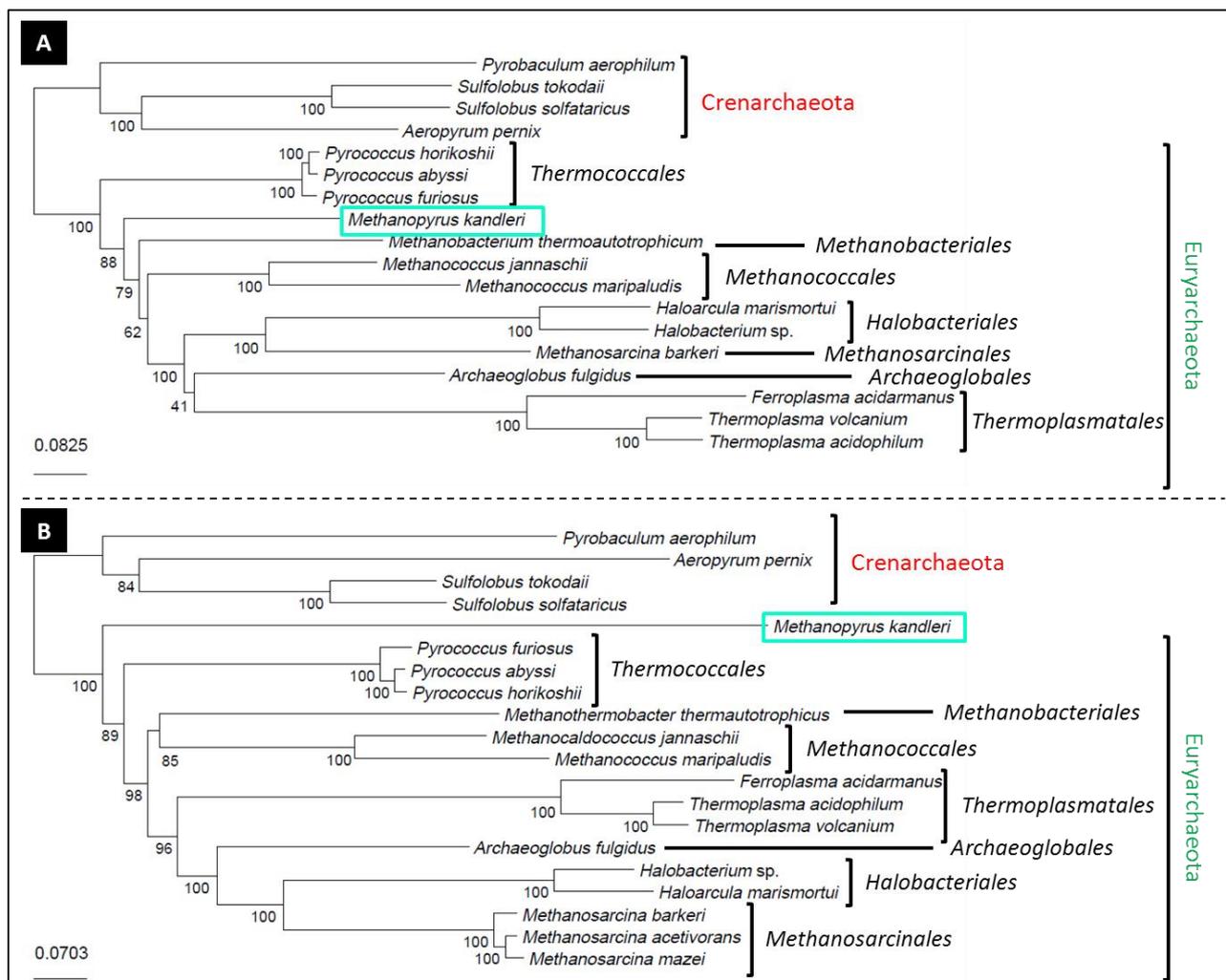


Figure 22. Position phylogénétique de *M. kandleri*. Phylogénies en maximum de vraisemblance adaptées de Brochier et collaborateurs 2004 (Brochier, Forterre, and Gribaldo 2004), racinées arbitrairement entre Euryarchaeota et Crenarchaeota. *M. kandleri* est entourée en bleu. A : Phylogénie inférée à partir de la concaténation de 53 protéines ribosomiques. *M. kandleri* branche après la divergence des Thermococcales au sein des Euryarchaeota. B : Phylogénie inférée à partir de la concaténation des sous-unités de l'ARN polymérase et de trois facteurs de transcription. *M. kandleri* branche entre les Crenarchaeota et les Euryarchaeota.

La position de *M. kandleri* est différente selon le jeu de données utilisé, basale sur une très longue branche avec les protéines impliquées dans la transcription ou au sein des Euryarchaeota avec les protéines ribosomiques. L'analyse des phylogénies individuelles des protéines impliquées dans la transcription montre que certaines sous-unités de l'ARN polymérase, A' et A'' par exemple, ont un signal fort pour placer *M. kandleri* à la base des Euryarchaeota sur une très longue branche, alors que la sous-unité B la place avec un fort support avec les *Methanococcales* et les *Methanobacteriales*, après la divergence des *Thermococcales*. Cette différence de position serait probablement le résultat d'un biais d'attraction de longues branches (LBA) dû à un taux d'évolution rapide des sous-unités de l'ARN polymérase. Par ailleurs, l'analyse d'un caractère rare appuie

l'hypothèse d'une divergence tardive de *M. kandleri* au sein des Euryarchaeota ; la sous-unité B est coupée en deux dans le génome de certaines euryarchées : les méthanogènes, les *Halobacteriales* et *Archaeoglobus fulgidus*. La coupure entre les deux parties du gène codant cette protéine dans ces génomes est exactement au même endroit de la séquence, il est donc plus parcimonieux de supposer qu'une seule coupure s'est produite chez l'ancêtre de ces organismes. Si *M. kandleri* se place à la base des Euryarchaeota, la coupure aurait dû se faire chez l'ancêtre de l'ensemble des Euryarchaeota, et les deux parties du gène auraient dû être rassemblées deux fois, chez les *Thermococcales* et chez les *Thermoplasmatales*. Si *M. kandleri* se place après la divergence des *Thermococcales*, alors une seule fusion est nécessaire pour expliquer cette répartition, rendant cette hypothèse plus parcimonieuse.

Un an plus tard, Bapteste et collaborateurs (Eric Bapteste, Brochier, and Boucher 2005) analysent l'évolution de la méthanogenèse et montrent que les protéines impliquées dans ce processus sont orthologues. La voie principale de méthanogenèse est la voie hydrogénotrophique, permettant la production de méthane à partir de CO₂ ou de formate, dans laquelle au moins 25 protéines sont impliquées. Certains des gènes codant ces protéines sont spécifiques aux méthanogènes alors que d'autres sont présents chez d'autres organismes. Néanmoins la phylogénie des gènes est congruente avec la phylogénie des organismes inférée sur les protéines ribosomiques, et montre deux groupes de méthanogènes au sein des Euryarchaeota. Le premier groupe est composé des *Methanopyrales* (dont *M. kandleri*), des *Methanococcales* et des *Methanobacteriales*. Les auteurs proposent de les rassembler sous le nom de « Méthanogènes Classe I » et ce groupe émergerait au sein des Euryarchaeota après la divergence des *Thermococcales*. Le deuxième groupe, rassemblant les *Methanosarcinales* et les *Methanomicrobiales*, serait les « Méthanogènes Classe II » et émergerait plus tardivement. La monophylie de chacune de ces classes n'est par contre pas clairement définie ou rejetée par ces analyses. Par ailleurs, il est notable que certains gènes impliqués dans la méthanogenèse sont présents chez *Archaeoglobus fulgidus* et chez les *Halobacteriales*, appuyant ainsi fortement l'hypothèse d'une apparition unique de la méthanogenèse au sein des Euryarchaeota avec des pertes différentielles dans les groupes phylogénétiquement placés entre les deux classes de méthanogènes. L'apparition de la méthanogenèse serait donc relativement tardive dans l'histoire évolutive des archées. L'impact de la production biologique de méthane sur l'atmosphère terrestre est très significatif en tant que gaz avec un fort effet de serre (cf. Introduction B.1.a). Son apparition est donc importante en ce qui concerne l'évolution des conditions de vie sur Terre.

La position de Nanoarchaeum equitans : nouveau phylum ou artefact ?

Lors de sa découverte l'archée symbiote *N. equitans* avait été proposée comme étant représentante d'un nouveau phylum d'archées, les Nanoarchaeota (H. Huber et al. 2002). Le séquençage de son génome est publié en 2003 par Waters et collaborateurs (Waters et al. 2003) et l'analyse phylogénétique de la concaténation de 35 protéines ribosomiques en utilisant les eucaryotes comme groupe extérieur montre sa position basale chez les archées. Brochier et collaborateurs réanalysent la position de cette espèce en 2005 en utilisant différentes méthodes (Brochier, Gribaldo, Zivanovic, Confalonieri, & Forterre, 2005). D'une part, différentes concaténations de protéines ribosomiques sont construites : une première contenant toutes les protéines ribosomiques, la même moins 9 protéines dont les phylogénies individuelles ont montré de possibles transferts horizontaux depuis les Crenarchaeota vers *N. equitans* (sa symbiose avec la crenarchée *I. hospitalis* pourrait expliquer facilement ces transferts) et enfin deux concaténations composées respectivement des protéines de la petite et de la grande sous-unités du ribosome. La première concaténation et la concaténation des protéines de la grande sous-unité du ribosome ont la topologie attendue pour les archées et *N. equitans* émerge entre les Euryarchaeota et les Crenarchaeota. La concaténation des protéines de la petite sous-unité du ribosome place par contre *N. equitans* au sein des Euryarchaeota comme groupe frère des *Thermococcales*, le reste de la phylogénie est celle attendue. Le signal contradictoire entre les deux sous-unités du ribosome remet en question la pertinence de l'utilisation conjointe de toutes les protéines ribosomiques dans ce cas, et l'analyse des phylogénies individuelles de ces protéines permet d'en écarter neuf ; la concaténation résultante donne une phylogénie où *N. equitans* est groupe frère des *Thermococcales* au sein des Euryarchaeota avec un support aux branches assez fort (Figure 23).

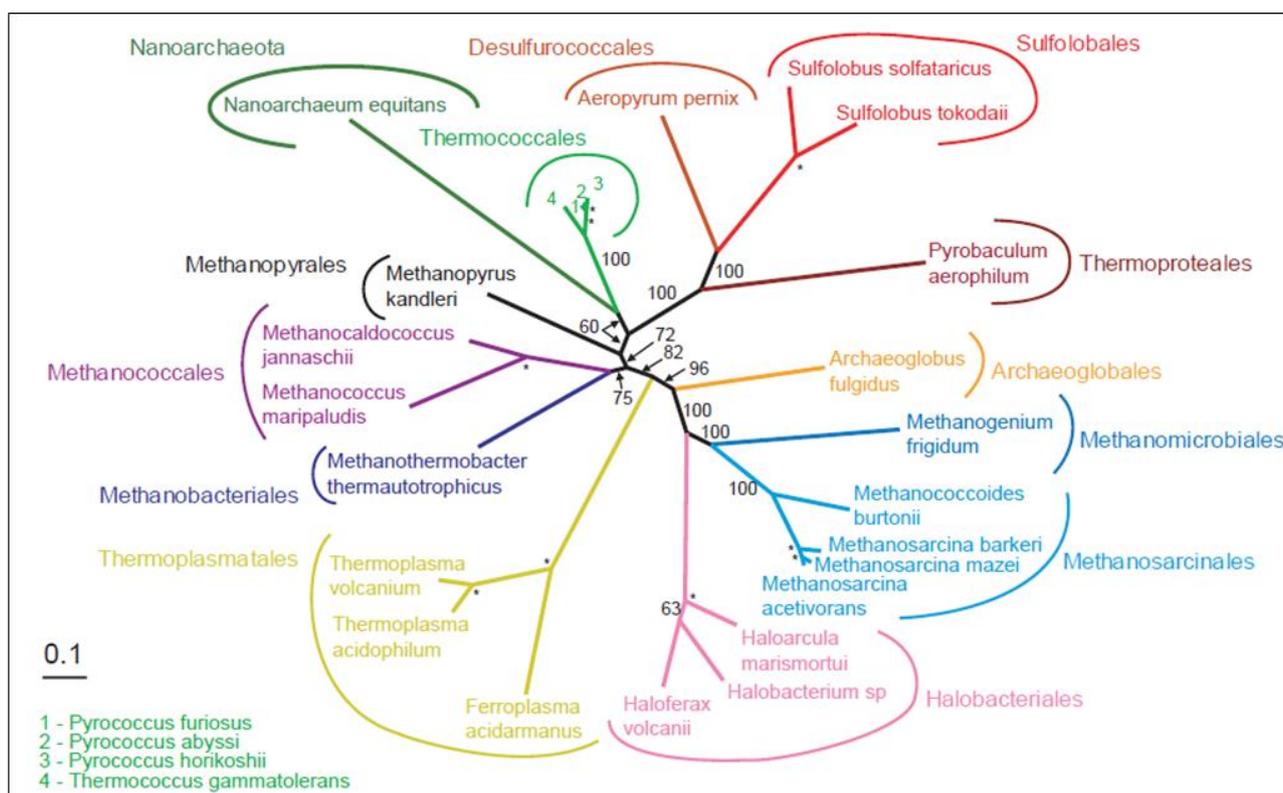


Figure 23. Nanoarchaeum equitans : une Euryarchaeota. Phylogénie de 41 protéines ribosomiques sélectionnées pour l'analyse de la position de *N. equitans* en maximum de vraisemblance par Brochier et collaborateurs 2005 (Brochier et al. 2005). On peut voir que *N. equitans* se place en groupe frère des Thermococcales avec un support de 60% et au sein des Euryarchaeota avec un support de 100%.

L'analyse de la phylogénie d'autres protéines, considérées comme des bons marqueurs de l'histoire évolutive des espèces, comme les facteurs d'élongation EF-1 α et EF2, ou celle de la sous-unité A de la topoisomérase VI et de la reverse-gyrase placent également *N. equitans* comme groupe frère des *Thermococcales*. De même, le premier résultat obtenu par BLAST à partir de l'ensemble de ses gènes est d'origine thermococcale dans 25 % des cas environ. L'ensemble de ces résultats semble bien confirmer l'appartenance de *N. equitans* aux Euryarchaeota, sa position basale sur une longue branche étant très probablement le résultat d'une LBA avec le groupe extérieur eucaryote provoquée par une évolution très rapide classique des symbiotes (Philippe and Laurent 1998).

La position basale de *N. equitans* chez les archées avait naturellement conduit certains auteurs à proposer que *N. equitans* soit représentante d'une ancienne lignée d'archées possédant ses caractéristiques génomiques particulières (H. Huber et al. 2002). Néanmoins, ces caractéristiques peuvent être interprétées comme le résultat d'une évolution rapide de leur génome par réduction, déjà observées chez des bactéries symbiotiques (Waters et al. 2003; Boucher, Doolittle, and Raynes 2002). Malgré cela, DiGiulio a proposé que ces caractéristiques de *N. equitans* soient en fait des caractères ancestraux au moins aux archées (Di Giulio 2008a; Di Giulio 2008b) et peut être même à l'ensemble du vivant (Di Giulio 2007). Ces caractères que DiGiulio considère comme ancestraux

sont les gènes codant les ARNt coupés en deux et l'absence d'opérons. Ils feraient de *N. equitans* un « génome fossile » (Di Giulio 2006). Les opérons et l'unification des gènes codant les ARNt seraient des caractères complexes apparus dans un second temps. Cette hypothèse ne repose pas sur des analyses phylogénétiques. Les multiples systèmes ayant la même histoire évolutive analysés par Brochier et collaborateurs corroborent fortement l'hypothèse opposée selon laquelle *N. equitans* est une euryarchée probablement groupe frère des *Thermococcales* et que ses caractéristiques génomiques sont le résultat d'une évolution symbiotique. La récente analyse du génome de « Nanoarchaeote Nst1 » (Podar et al. 2013) montre une évolution par réduction de génome du même type que *N. equitans*, mais beaucoup moins avancée, ce qui apporte une preuve de plus de leur évolution atypique.

Nouveaux génomes et nouveaux marqueurs

En 2005, 25 génomes d'archées sont disponibles, 4 Crenarchaeota et 21 Euryarchaeota (incluant *N. equitans*), représentant 13 ordres. Les trois précédentes études portant sur la phylogénie des archées décrites plus haut portaient sur des points précis de cette histoire évolutive (position de *M. kandleri*, évolution de la méthanogenèse et position de *N. equitans*) ; le nombre de génomes d'archées disponibles depuis la dernière étude de leur phylogénie en général (Matte-Tailliez et al. 2002) a presque doublé à ce moment-là, et Brochier et collaborateurs font alors une mise à jour des jeux de données des protéines ribosomiques et des sous-unité de l'ARN polymérase pour tester la position des nouveaux génomes disponibles et vérifier la robustesse des topologies obtenues jusqu'ici (Brochier, Forterre, and Gribaldo 2005). Deux concaténations sont construites, la première nommée « traduction » à partir de 53 protéines ribosomiques, et la seconde « transcription » à partir des 15 protéines impliquées dans la transcription (12 sous-unité de l'ARN polymérase ; A', A'', B, D, E', E'', F, H, K, L, N, P ; et les facteurs de transcription NusA, NuG et TFS). Les deux phylogénies obtenues sont congruentes, excepté en ce qui concerne la position de *M. kandleri*, qui se place à la base des Euryarchaeota dans l'arbre de « transcription » et après les *Thermococcales* dans l'arbre de « traduction », comme attendu. La congruence de ces deux phylogénies tendrait à prouver qu'un ensemble de gènes, le « core génome », évolue principalement par héritage vertical chez les archées et enregistre l'histoire évolutive de ces organismes. Cette histoire serait donc reconstructible malgré les transferts horizontaux de gènes existants chez les archées. Néanmoins, ces deux systèmes, le ribosome et l'ARN polymérase sont impliqués dans des processus informationnels de la cellule (transcription et traduction) et qui, de plus, sont couplés (French et al. 2007). Il n'est pas exclu que l'histoire évolutive qu'ils reflètent soit l'histoire des systèmes informatifs mais que d'autres systèmes, métaboliques ou liés à d'autres processus cellulaires,

n'aient pas eu la même évolution.

L'année suivante, Gribaldo et Brochier-Armanet font une revue des travaux sur la phylogénie des archées (Simonetta Gribaldo and Brochier-Armanet 2006), et actualisent les jeux de données transcription et traduction avec quatre nouveaux génomes (Figure 24) et étudient la phylogénie d'autres systèmes, le complexe SRP (Signal Recognition Particule) et l'exosome.

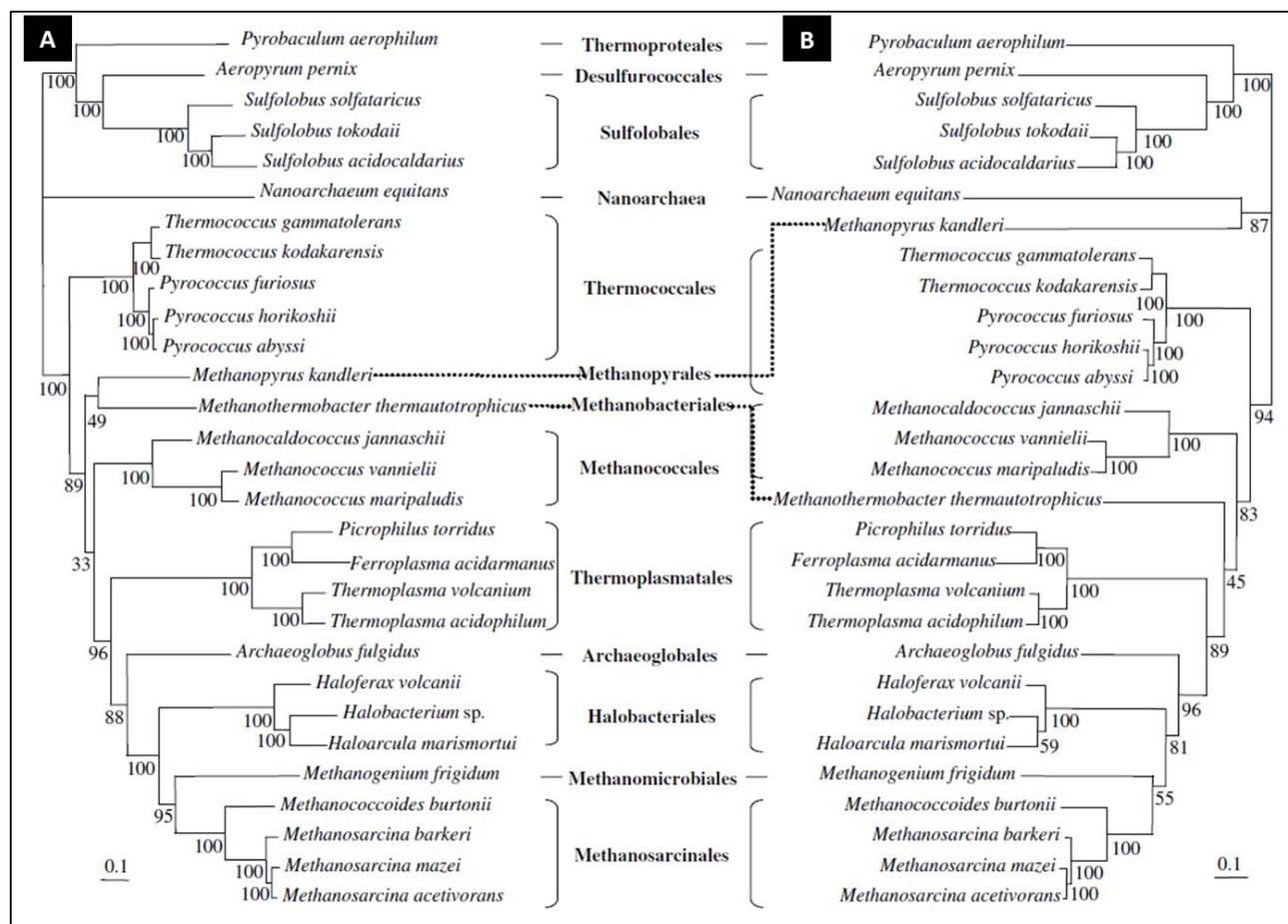


Figure 24. Phylogénie des archées en 2006. Phylogénies en maximum de vraisemblance non racinées de Gribaldo et Brochier-Armanet en 2006 (Simonetta Gribaldo and Brochier-Armanet 2006). A : Phylogénie inférée à partir de la concaténation de 53 protéines ribosomiques. B : Phylogénie inférée à partir de la concaténation des sous-unités de l'ARN polymérase et de trois facteurs de transcription.

Les deux phylogénies de la Figure 24 sont congruentes sur leurs topologies globales et très proches de celles obtenues lors des études précédentes, mais mieux soutenues. On retrouve les Euryarchaeota et les Crenarchaeota, et *N. equitans* se place entre les deux. Au sein des Crenarchaeota, *Aeropyrum pernix*, représentant des *Desulfurococcales*, se place en groupe frère des *Sulfolobales*, et *Pyrobaculum aerophilum*, représentant les *Thermoproteales* est le premier à diverger. Ces relations de parenté ont déjà été observées dans les études précédentes. En ce qui concerne les Euryarchaeota, les ordres représentés par plusieurs organismes sont monophylétiques :

les *Thermococcales*, les *Methanococcales*, les *Thermoplasmatales*, les *Halobacteriales*, et les *Methanosarcinales*. L'ordre de divergence chez les euryarchées est aussi très ressemblant entre ces deux phylogénies. Les méthanogènes classe II (*Methanomicrobiales* et *Methanosarcinales*) forment un groupe monophylétique et sont groupe frère des *Halobacteriales*. Ce sous ensemble est à son tour groupe frère d'*Archaeoglobus fulgidus*, puis des *Thermoplasmatales*. On trouve par contre des différences en ce qui concerne les relations entre les méthanogènes classe I (*Methanopyrales*, *Methanococcales* et *Methanobacteriales*). Dans la phylogénie « traduction », *M. kandleri* et *Methanothermobacter thermoautotrophicus* (une méthanobactériale) sont groupes frères et émergent juste après les *Thermococcales*, les *Methanococcales* se plaçant entre le groupe formé par *M. kandleri* et *M. thermoautotrophicus* et le reste des euryarchées. Dans la phylogénie « transcription », *M. kandleri* est groupe frère de *N. equitans*, et *M. thermoautotrophicus* se place entre les *Methanococcales* et le reste des euryarchées. La position de *M. kandleri* à la base des euryarchées est très probablement le résultat d'une LBA comme expliqué plus tôt dans l'arbre « transcription », mais de fait, les relations entre les méthanogènes classe I ne sont pas résolues. Il est aussi possible qu'un manque de signal phylogénétique provoque cette irrésolution comme le suggèrent les faibles soutiens de ces nœuds.

Le système SRP et l'exosome sont deux complexes protéiques impliqués dans d'autres mécanismes cellulaires que les processus informationnels. Ils sont présents chez toutes les archées et composés de plusieurs protéines ; dans chacun de ces deux complexes, les protéines partagent très probablement une même histoire évolutive puisque agissant en complexe dans la cellule. Il était donc intéressant de les comparer avec les systèmes informationnels de transcription et traduction. Le système SRP est composé de deux protéines paralogues, Srp54 et Srp19, d'un récepteur FtsY, et d'un ARN 7S. Srp54 et FtsY sont présents chez toutes les archées à l'exception de *N. equitans*, alors que Srp19 est absent dans plusieurs ordres (*Thermococcales*, *Thermoproteales*, et *N. equitans*) : ce dernier sera donc écarté de l'analyse. Les phylogénies individuelles de Srp54 et de FtsY ont des topologies proches de la phylogénie « informationnelle », avec une dichotomie Crenarchaeota/Euryarchaeota, la monophylie des ordres, et le même ordre de divergence chez les Crenarchaeota. Par contre, les relations entre euryarchées sont très mal résolues. Aucun transfert horizontal n'est clairement détecté, mais le manque de signal dû au peu de positions des alignements individuels est probablement la cause de cette irrésolution. La concaténation des deux protéines montre une phylogénie très proche de la phylogénie « informationnelle », avec quelques points de divergences : les *Methanobacteriales* et *Methanococcales* sont groupes frères entre eux, et leur groupe est frère des *Thermococcales*, alors que *M. kandleri* branche à la base des euryarchées.

L'exosome est un complexe de dégradation des ARN dans la cellule. Chez les archées il est

formé de quatre protéines : Rrp41 et Rrp42, qui composent l'anneau central ayant une activité exonucléase, et Rrp4p et Cs14, deux protéines périphériques. Il est absent chez les *Halobacteriales* et chez les *Methanococcales*, et Cs14 est absent chez *N. equitans*. Les phylogénies individuelles ayant permis de détecter un transfert horizontal des crenarchées vers *N. equitans*, la séquence en question a donc été retirée. La phylogénie obtenue à partir de la concaténation des quatre protéines soutient les mêmes relations que pour les protéines du complexe SRP, avec un faible support et *N. equitans*, absente dans la phylogénie « SRP », branche à la base des archées. Cette faible résolution est probablement le résultat d'un manque de signal et des données manquantes pour Cs14 en ce qui concerne la position de *N. equitans*.

Les phylogénies de ces deux systèmes, bien que faiblement résolues et incomplètes pour l'exosome, montrent un signal très proche de celui obtenu à partir des protéines informationnelles. L'ajout de ces nouveaux marqueurs ne permet pas d'apporter plus d'éléments en ce qui concerne les relations problématiques entre certains groupes d'archées, mais, et c'est le plus important, confirme que l'histoire évolutive des archées est inscrite dans un certain nombre de protéines héritées verticalement.

Nouveaux phyla et phylogénie de référence actuelle

En 2008, le phylum des Thaumarchaeota est proposé par Brochier-Armanet et collaborateurs (Brochier-Armanet et al. 2008) et depuis plusieurs génomes de Thaumarchaeota ont été publiés. La même année, le génome de '*Ca. Korarchaeum cryptophilum*' est publié (Elkins et al. 2008). Enfin, en 2011, le génome de '*Ca. Caldiarchaeum subterraneum*' est publié (Nunoura et al. 2011) et le phylum Aigarchaeota est proposé par ses auteurs. Le séquençage de génomes représentant des lignées basales de la phylogénie des archées est donc très important pour comprendre les relations anciennes entre les différents phyla et mieux connaître la nature du dernier ancêtre commun aux archées.

Comme déjà expliqué précédemment, les Thaumarchaeota correspondent au groupe I d'archées mésophiles proches des crenarchées, décrites indépendamment par DeLong et Furhmann et collaborateurs en 1992. C'est un groupe d'organismes nombreux, en termes de biomasse et de diversité. Ils partagent la capacité d'oxydation anaérobie de l'ammonium, et leurs relations phylogénétiques avec les Crenarchaeota ne sont pas clairement résolues. La publication du génome de *Cenarchaeum symbiosum* (Hallam et al. 2006) permet enfin d'intégrer le groupe I à la phylogénie des archées en 2008 (Brochier-Armanet et al. 2008). L'arbre basé sur la concaténation des ARNr place le groupe I comme groupe frère des Crenarchaeota, et la racine de l'arbre est placée entre ce groupe et les Euryarchaeota. La phylogénie inférée à partir des protéines ribosomiques

(Figure 25) utilisant les eucaryotes comme groupe extérieur donne un résultat différent car *Cenarchaeum symbiosum* se place à la base des archées, avant l'émergence des Euryarchaeota et des Crenarchaeota.

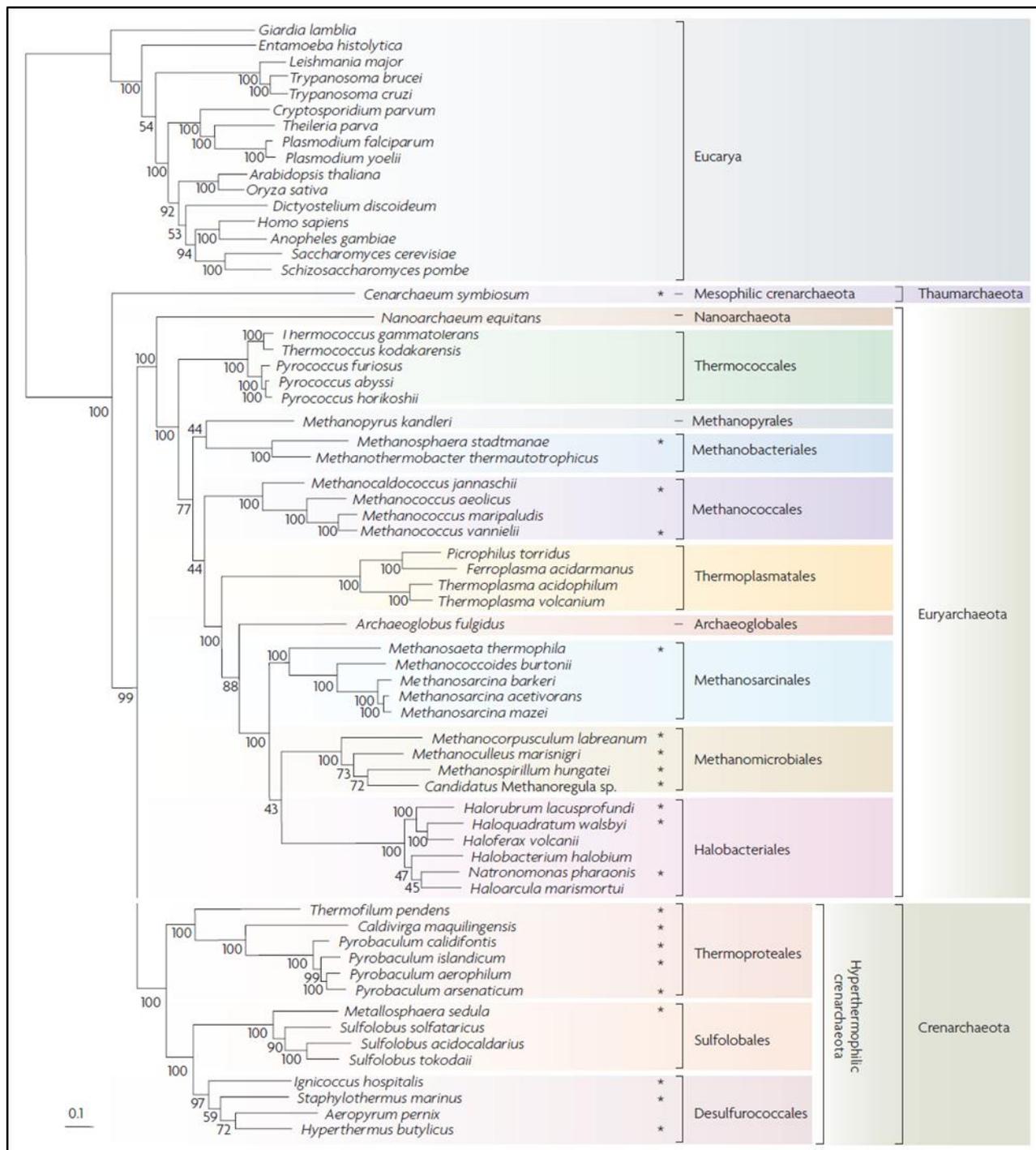


Figure 25. Les Thaumarchaeota : un nouveau phylum? Phylogénie en maximum de vraisemblance des protéines ribosomiques publiée par Brochier-Armanet et collaborateurs en 2008 (Brochier-Armanet et al. 2008). *Cenarchaeum symbiosum*, représentant les Thaumarchaeota, branche basalement au sein des archées.

Parallèlement à cette analyse phylogénétique, les auteurs ont fait une analyse du contenu en gènes dans le génome de *C. symbiosum*. L'analyse des COG, « Clusters of Orthologous Groups » (R L Tatusov et al. 2001) (qui permettent une classification des gènes orthologues selon leur fonction cellulaire) présents chez *C. symbiosum* montre qu'il partage un nombre de catégories fonctionnelles du même ordre avec les euryarchées (sept) et les crenarchées (un) que ces deux phyla entre eux (25), reflétant une grande différence biologique entre ces trois groupes. D'autres caractéristiques particulières placent ce génome à part des deux phyla d'archées : l'absence ou la présence de protéines ribosomiques particulières, l'absence de l'ADN topoisomérase IA (présente chez les autres archées, les bactéries et les eucaryotes) ou la présence de la ADN topoisomérase IB (absente chez les autres archées mais présente chez les eucaryotes). L'ensemble de ces caractéristiques, qui ne les rapprochent ni des Crenarchaeota ni des Euryarchaeota, et leurs particularités métaboliques sont donc des indices forts pour les considérer comme un groupe bien particulier. Ils ont conduit Brochier-Armanet et collaborateurs à proposer la création du nouveau phylum des Thaumarchaeota. La place de ce phylum au sein des archées n'est par contre pas résolue, puisque selon le groupe extérieur choisi, bactéries ou eucaryotes, la racine se place soit dans la branche menant aux Euryarchaeota, soit dans celle menant aux Thaumarchaeota, respectivement.

Du point de vue de la phylogénie des archées, la [Figure 25](#) montre un arbre raciné avec les séquences eucaryotes. Cet arbre est très intéressant puisque qu'il montre les relations phylogénétiques entre 48 archées très diverses. La monophylie des Crenarchaeota et des Euryarchaeota est retrouvée, mais aussi celle des différents ordres d'archées. En ce qui concerne les points problématiques abordés plus tôt, il est intéressant de constater que *N. equitans* se place à la base des Euryarchaeota avec un support maximal. La phylogénie inférée à partir des ARNr et racinée par les séquences bactériennes place *N. equitans* au même endroit, si ce n'est qu'elle groupe avec *M. kandleri*. Par contre, dans la phylogénie inférée à partir des protéines ribosomiques, *M. kandleri* branche avec les *Methanobacteriales* et les méthanogènes classe I ne sont toujours pas monophylétiques. De même les méthanogènes classe II, qui semblaient jusque là être monophylétiques, ne le sont pas, et les *Methanomicrobiales* sont groupes frères avec les *Halobacteriales*.

J'ai déjà décrit précédemment comment avaient été proposés les phyla Korarchaeota (Barns et al. 1996) et Aigarchaeota (Nunoura et al. 2011). En 2008 et 2011 sont publiées les séquences des génomes des deux premiers représentants de ces phyla, '*Ca. Korarchaeum cryptophylum*' et '*Ca. Caldiarchaeum subterraneum*', respectivement. En 2011, Brochier-Armanet et collaborateurs (Brochier-Armanet, Forterre, & Gribaldo, 2011) intègrent ces nouveaux génomes à la phylogénie

des archées, ainsi que de nouvelles espèces découvertes entre temps, telles que les ARMAN, le nouvel ordre de méthanogènes, les *Methanocellales*, *Acidilobus saccharovorans* (représentant des *Acidilobales*), *Aciduliprofundum boonei* (représentant du groupe d'euryarchées DHVE2) et un fragment du génome de l'organisme « uncultured marine DeepAnt-JyKC7 » (représentant non cultivé du groupe II d'Euryarchaeota) ainsi que les nouveaux génomes séquencés appartenant aux ordres déjà représentés précédemment. Une phylogénie est inférée à partir de 57 protéines ribosomiques présentes dans au moins 89 des 99 génomes d'archées utilisés (Figure 26). Les six phyla d'archées (y compris le phylum Nanoarchaeota, bien qu'il ait été clairement remis en question) y sont représentés et le nombre de génomes utilisés par rapport à la phylogénie précédente a doublé. Néanmoins cette phylogénie n'est pas racinée, et il est donc impossible de définir les relations anciennes entre phyla.

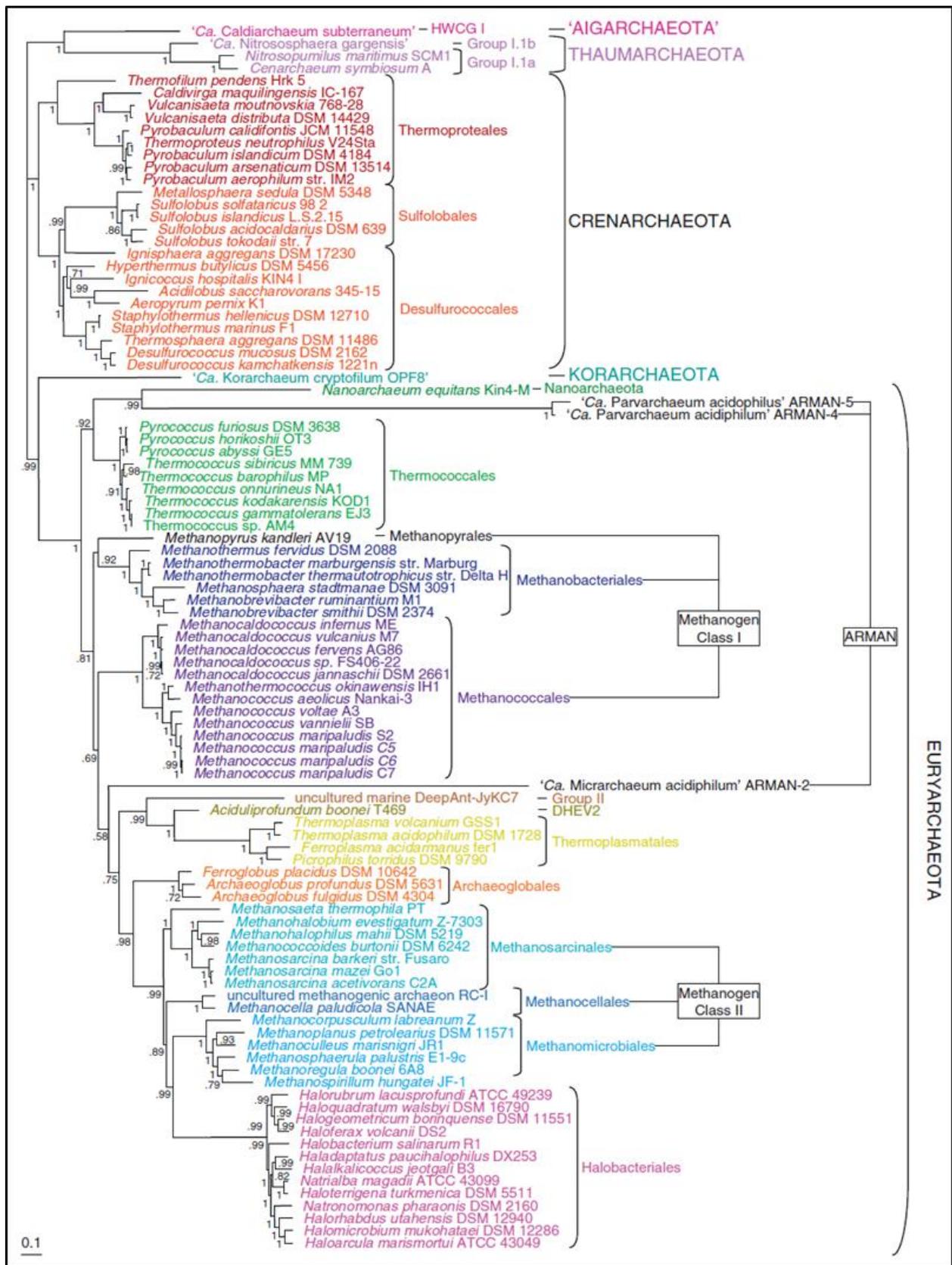


Figure 26. Phylogénie des archées en 2011 (référence actuelle). Phylogénie non racinée en inférence bayésienne publiée par Brochier-Armanet et collaborateurs en 2011 (Brochier-Armanet, Forterre, and Gribaldo 2011), inférée à partir d'un alignement de 57 protéines ribosomiques de 99 espèces (5838 positions).

La nouvelle phylogénie vient confirmer un certain nombre de points déjà observés dans les phylogénies précédentes : la monophylie des Euryarchaeota, des Crenarchaeota et des Thaumarchaeota ; la monophylie des différents ordres d'archées ; la position de *N. equitans* en groupe frère des *Thermococcales* et la position de *M. kandleri* au sein des Euryarchaeota après la divergence des *Thermococcales*. En ce qui concerne les relations entre phyla (excepté les Nanoarchaeota), on peut noter la relation entre les Thaumarchaeota et les Aigarchaeota ('*Ca. Caldiarchaeum subterraneum*'), avec un support maximal. Hormis cette relation privilégiée, les relations entre phyla ne peuvent pas être inférées clairement à partir de cet arbre étant donné qu'il n'est pas raciné. Il est donc impossible de savoir quel est l'ordre de divergence entre phyla, ou si un regroupement existe entre différents phyla.

Cependant, elle permet d'apporter de nouveaux éléments quant à la position de certaines espèces. Par exemple *Acidilobus saccharovorans*, une crenarchée appartenant à l'ordre récemment proposé *Acidilobales* (Prokofeva et al. 2009), se trouve en fait à l'intérieur des *Desulfurococcales* à proximité d'*Aeropyrum pernix*. Le groupe des *Acidilobales* semble donc faire partie intégrante des *Desulfurococcales* et ne pas être un ordre à part entière. L'ordre proposé récemment des *Fervidicoccales* (Perevalova et al. 2010) serait phylogénétiquement proche des *Acidilobales*, il est fort probable que sa situation soit donc la même. La position d'*Ignisphaera aggregans* est, par contre, problématique car elle se place à la base des *Sulfolobales* alors qu'elle est décrite comme appartenant aux *Desulfurococcales*, faisant de cet ordre un groupe paraphylétique contenant les *Sulfolobales*.

Chez les Euryarchaeota, en plus du placement de *N. equitans* avec les *Thermococcales*, cette phylogénie permet de placer le nouveau génome d'*Aciduliprofundum boonei*, du groupe DHVE2, en groupe frère des *Thermoplasmatales*, l'ensemble formant un groupe monophylétique avec le génome « uncultured marine DeepAnt-JyKC7 » (représentant non cultivé du groupe II d'Euryarchaeota). Étant donné la large diversité d'organismes non cultivés connus pour le groupe II et le groupe DHVE2 (Christa Schleper, Jurgens, and Jonuscheit 2005), il est attendu que cette région de l'arbre des Euryarchaeota accueille prochainement des nouveaux représentants. Par contre, la position des génomes des ARMAN, '*Ca. Micrarchaeum acidiphilum* ARMAN-2', '*Ca. Parvarchaeum acidiphilum* ARMAN-4' et '*Ca. Parvarchaeum acidophilus* ARMAN-5', n'est pas résolue. D'une part, les deux '*Ca. Parvarchaeum*' se placent ensemble, ce qui est attendu étant donné la forte similarité de leurs ARNr SSU, mais sur une très longue branche, à côté de *N. equitans*. Comme cela a déjà été abordé, ces organismes de petite taille ont un génome réduit (de l'ordre de 1 mégabase) et il est possible qu'ils aient eu une évolution rapide liée à un mode de vie

symbiotique ou parasitique, entraînant leur position sur une longue branche très instable. La position des '*Ca. Parvarchaeum*' à côté de *N. equitans*, elle-même sur une longue branche, semble donc artificielle. En ce qui concerne la troisième ARMAN, '*Ca. Micrarchaeum acidiphilum*', elle se place après la divergence des *Methanococcales* et juste avant la divergence du groupe composé des *Thermoplasmatales*, du groupe DHVE2 et du groupe II, mais cette position n'est pas soutenue et devra être clarifiée. Enfin, ni les méthanogènes classe I ni les classe II ne sont monophylétiques, comme dans la dernière phylogénie des archées proposée en 2008 (Brochier-Armanet et al., 2008), mais contrairement à d'autres phylogénies plus anciennes (Simonetta Gribaldo and Brochier-Armanet 2006). Par ailleurs, le nouvel ordre d'archées méthanogènes *Methanocellales*, représenté par deux espèces, se place entre les *Methanosarcinales* et les *Methanomicrobiales*, elles mêmes groupe frère des *Halobacteriales*.

Actuellement, cette phylogénie des archées peut être considérée comme phylogénie de référence des espèces, dans la mesure où elle est centrée sur les archées, inférée à partir de plusieurs marqueurs fiables, assez bien soutenue et représentative de la diversité de ce domaine.

Le Tableau 1 récapitule les études de la phylogénie des archées que je viens de citer, avec le nombre de génomes utilisés par phyla d'archées, le nombre de protéines ribosomiques et d'autres protéines utilisées.

Article	Génomes	Protéines ribosomiques	Autres marqueurs protéiques
Matte-Taille et collaborateurs 2002 (Matte-Tailliez et al. 2002)	3 crenarchées 10 euryarchées	53	ARN polymérase : Sous-unité RpoB
Brochier et collaborateurs 2004 (Brochier et al., 2004)	4 crenarchées 14 euryarchées	53	ARN polymérase : 12 sous-unités Facteurs de transcription : 3
Baptiste et collaborateurs 2005 (Eric Baptiste, Brochier, and Boucher 2005)	4 crenarchées 19 euryarchées (<i>⊂ N. equitans</i>)	53	Protéines impliquées dans la méthanogenèse
Brochier et collaborateurs 2005 (Brochier et al. 2005)	4 crenarchées 21 euryarchées (<i>⊂ N. equitans</i>)	50	Facteurs d'élongation : EF-1 α et EF-2 Topoisomérase VI : Sous-unité A Reverse gyrase
Brochier et collaborateurs 2005 (Brochier, Forterre, and Gribaldo 2005)	4 crenarchées 21 euryarchées (<i>⊂ N. equitans</i>)	53	ARN polymérase : 12 sous-unités Facteurs de transcription : 3
Gribaldo & Brochier-Armanet 2006 (Simonetta Gribaldo and Brochier-Armanet 2006)	5 crenarchées 24 euryarchées (<i>⊂ N. equitans</i>)	53	ARN polymérase : 12 sous-unités Facteurs de transcription : 3 Système SRP : Srp 54 et FtsY Exosome : 4 sous-unités
Brochier-Armanet et collaborateurs 2008 (Brochier-Armanet et al. 2008)	14 crenarchées 1 thaumarchée 33 euryarchées (<i>⊂ N. equitans</i>)	53	
Brochier-Armanet et collaborateurs 2011 (Brochier-Armanet, Forterre, and Gribaldo 2011)	24 crenarchées 3 thaumarchées 1 aigarchée 1 korarchée 70 euryarchées (<i>⊂ N. equitans</i>)	57	

Tableau 1. Etudes de la phylogénie des archées : génomes et marqueurs utilisés.

b. La racine de l'arbre des archées et la relation avec les autres domaines

Comme nous venons de le voir, l'étude de la phylogénie des archées est très souvent passée par l'utilisation des protéines informationnelles, particulièrement des protéines ribosomiques et sans groupe extérieur. Les relations anciennes chez les archées, en l'occurrence entre les différentes grandes lignées d'archées et les premières divergences au sein de ces groupes, sont importantes pour comprendre, d'une part, la nature du dernier ancêtre commun des archées (LACA pour « Last Archaeal Common Ancestor ») et, d'autre part, les relations entre ce domaine du vivant et les deux autres domaines, les bactéries et les eucaryotes.

En ce qui concerne les relations entre les archées et les autres domaines du vivant, les archées partagent avec les eucaryotes un certain nombre de caractéristiques les rapprochant,

particulièrement les systèmes « informationnels », à savoir ceux impliqués dans la transmission et l'expression de l'information dans la cellule (réplication, transcription et traduction), qui sont très similaires en termes de composition protéique et de similarité de séquences entre archées et eucaryotes par rapport aux bactéries. Leur relation privilégiée a été établie par le positionnement de la racine de l'arbre du vivant dans la branche menant aux bactéries dès 1989, grâce aux couples de paralogues universels, EF-1 α et EF-2, les sous-unités α et β de l'ATPase F₁ (Iwabe et al. 1989; Gogarten et al. 1989), et Srp54/Ffh et SR α /FtsY (S Gribaldo and Cammarano 1998). Même si la relation phylogénétique entre archées et eucaryotes semble acquise, la nature de cette relation n'est pas pour autant résolue. Deux écoles existent : la première propose que les archées et les eucaryotes soient groupes frères, comme l'illustre l'arbre du vivant de Woese en 1990 (Figure 4). Il y aurait ainsi trois domaines primaires : c'est l'hypothèse « 3D ». La seconde propose que les eucaryotes soient le résultat de l'association d'une archée, qui aurait donné la majeure partie des gènes informationnels eucaryotes, les plaçant ainsi phylogénétiquement proches des archées, et d'une bactérie. Il n'y aurait que deux domaines primaires : c'est l'hypothèse « 2D ». De nombreuses études phylogénomiques ont tenté de résoudre les relations entre les trois domaines du vivant comme l'illustre le Tableau 2. La plupart de ces études ont été faites à partir d'un nombre réduit de marqueurs, majoritairement des protéines ribosomiques, et en prenant en compte l'ensemble du vivant, archées, eucaryotes et bactéries. Les résultats varient selon les études, et le débat n'est pas tranché.

Article	Génomes	Marqueurs	Hypothèse soutenue
Harris et collaborateurs 2003 (Harris et al. 2003)	25 bactéries 1 crenarchée 7 euryarchées 3 eucaryotes	50 (\approx 29 protéines ribosomiques)	3D
Rivera & Lake 2004 (Rivera and Lake 2004)	2 bactéries 1 crenarchée 2 euryarchées 2 eucaryotes	Génome complet d' <i>Archaeoglobus fulgidus</i>	2D (Eucaryotes groupe frère des Crenarchaeota)
Ciccarelli et collaborateurs 2006 (Ciccarelli et al. 2006)	150 bactéries 4 crenarchées 14 euryarchées 23 eucaryotes	31 (\approx 23 protéines ribosomiques)	3D
Yutin et collaborateurs 2008 (Yutin et al. 2008)	Variable selon les marqueurs	136 (\approx 31 protéines ribosomiques)	3D
Pisani et collaborateurs 2007 (Pisani, Cotton, and McInerney 2007)	97 bactéries 4 crenarchées 17 euryarchées 2 eucaryotes	Données non disponibles	2D (Eucaryotes groupe frère des Thermoplasmatales)
Cox et collaborateurs 2008 (Cox et al. 2008)	10 bactéries 3 crenarchées 11 euryarchées 16 eucaryotes	45 (\approx 20 protéines ribosomiques, 2 ARN ribosomiques)	2D (Eucaryotes groupe frère des Crenarchaeota)
Foster et collaborateurs 2009 (Foster, Cox, and Embley 2009)	8 bactéries 8 crenarchées 2 thaumarchées 6 euryarchées 11 eucaryotes	41 (\approx 20 protéines ribosomiques)	2D (Eucaryotes groupe frère du groupe Crenarchaeota+Thaum archaeota)
Guy & Ettema 2010 (Guy and Ettema 2011)	10 bactéries 15 crenarchées 3 thaumarchées 1 aigarchée 1 korarchée 12 euryarchées (\approx <i>N. equitans</i>) 7 eucaryotes	26 (\approx 20 protéines ribosomiques)	2D (Eucaryotes groupe frère des « TACK »)
Williams et collaborateurs 2012 (Williams et al. 2012)	8 bactéries 8 crenarchées 3 thaumarchées 1 aigarchée 1 korarchée 6 euryarchées 11 eucaryotes	29 (\approx 13 protéines ribosomiques)	2D (Eucaryotes groupe frère des « TACK »)
Lasek-Nesselquist & Gogarten 2013 (Lasek-Nesselquist and Gogarten 2013)	73 bactéries 17 crenarchées 3 thaumarchées 1 korarchée 30 euryarchées (\approx <i>N. equitans</i>) 38 eucaryotes (variable selon les analyses)	85 protéines ribosomiques (variable selon les analyses)	Tendance 2D (variable selon les analyses)

Tableau 2. Analyses phylogénomiques sur l'arbre du vivant. Adapté de Gribaldo et collaborateurs 2010 (Simonetta Gribaldo et al. 2010).

Dans ce contexte, les relations anciennes chez les archées représentent un problème essentiel. En effet, il est important de connaître l'ordre de divergence des différents phyla d'archées pour pouvoir replacer les eucaryotes. Si nous sommes dans l'hypothèse « 2D », les eucaryotes brancheront au sein des archées, et il sera possible de savoir quelles sont les lignées d'archées les plus proches des eucaryotes ; dans une hypothèse « 3D », la position de la racine de l'arbre des archées permettra de savoir si les eucaryotes sont bien groupe frère de l'ensemble des archées ou s'ils émergent à l'intérieur de l'arbre des archées. Très peu d'études ont été faites ayant pour but de positionner la racine de l'arbre des archées. En 2008, Brochier-Armanet et collaborateurs (Brochier-Armanet et al. 2008) ont inféré la phylogénie des archées avec un groupe extérieur pour pouvoir comprendre les relations entre les Thaumarchaeota, les Crenarchaeota et les Euryarchaeota, mais un seul génome complet de Thaumarchaeota était disponible et le groupe extérieur choisi pour l'arbre inféré à partir des protéines ribosomiques était les eucaryotes, ce qui ne permet pas de résoudre les relations entre les deux domaines. Un arbre inféré à partir des ARNr avait pour groupe extérieur les bactéries, mais comme je l'ai déjà expliqué plus tôt, ces marqueurs sont insuffisants pour une bonne résolution de l'arbre des archées, et particulièrement des relations anciennes.

A contrario, des analyses à large échelle ont été faites pour essayer de résoudre les relations entre les trois domaines (Tableau 2). Les études les plus récentes prennent en compte les phyla d'archées pour lesquels des génomes ont été publiés durant ces trois dernières années. Par exemple, en 2011, Guy et Ettema (Guy and Ettema 2011) font une étude phylogénomique utilisant 26 protéines universelles (dont 20 protéines ribosomiques et deux sous-unités de l'ARN polymérase) (Tableau 2) sur 48 génomes (31 archées, sept eucaryotes et dix bactéries comme groupe extérieur). Ce jeu de données a été traité avec plusieurs méthodes, et le résultat place toujours les eucaryotes au sein des archées, soit en groupe frère du groupe monophylétique formé par les Thaumarchaeota et les Aigarchaeota, soit entre les Euryarchaeota (y compris *N. equitans*) et les autres phyla (Thaumarchaeota, Aigarchaeota, Crenarchaeota, Korarchaeota). Ils proposent ainsi de définir le groupe « TACK », acronyme des noms de ces quatre phyla, qui serait particulièrement relié aux eucaryotes. Néanmoins, ce travail pose un certain nombre de problèmes : il est basé sur un échantillonnage taxonomique très réduit, particulièrement en ce qui concerne les eucaryotes et les archées, principaux objets de l'étude ; *N. equitans* est intégré au jeu de données, alors que cette espèce est connue pour être soumise à l'artefact d'attraction de longues branches en raison de son évolution rapide ; les marqueurs utilisés sont très peu nombreux. Ce dernier point est particulièrement problématique. Les jeux de données utilisés ont été définis à partir des travaux antérieurs de Cox et collaborateurs (Cox et al. 2008) et de Ciccarelli et collaborateurs (Ciccarelli et al. 2006), qui avaient déjà été critiqués pour leur manque de fiabilité (Simonetta Gribaldo et al.

2010). De plus, l'utilisation de ces marqueurs sur l'ensemble du vivant réduit d'autant plus le nombre de positions utilisables pour l'inférence phylogénétique. Enfin, le parti pris des auteurs est clairement affiché dès le résumé de l'article et se place dans l'hypothèse « 2D ».

Dans l'article de Rinke et collaborateurs publié très récemment (Rinke et al. 2013), les auteurs présentent le séquençage en « single-cell » de 201 archées et bactéries non cultivées provenant de divers environnements. En termes d'apport de nouvelles données, cet article est extrêmement novateur car les organismes séquencés ont été choisis justement pour leur appartenance à des lignées peu connues, ce qu'ils appellent la « matière noire microbienne ». Les auteurs ont de plus utilisé ces nouveaux génomes complets pour étudier la phylogénie des archées et des bactéries. Une première phylogénie a été inférée à partir de l'ARNr SSU de l'ensemble des organismes séquencés dans ce travail et d'organismes de référence déjà disponibles. Elle permet aux auteurs de raciner des archées avec les bactéries et, réciproquement, les bactéries avec les archées. La racine de l'arbre des archées se place entre un groupe composé des phyla Thaumarchaeota, Aigarchaeota, Crenarchaeota et Korarchaeota (« TACK ») augmenté d'un certain nombre de lignées nouvelles, et un groupe composé des Euryarchaeota, comprenant aussi un certain nombre de nouvelles lignées. Ensuite, une phylogénie a été inférée pour chaque domaine à partir d'une concaténation de 38 protéines utilisées régulièrement pour ce type d'analyse (dont 27 protéines ribosomiques). En ce qui concerne les archées, 220 génomes ont été utilisés. Les auteurs retrouvent la monophylie du groupe « TACK » proposé par Guy et Ettema, qui se place en groupe frère des Euryarchaeota. Un autre groupe composé de plusieurs lignées anciennes et nouvelles semble avoir divergé avant les Euryarchaeota et les « TACK ». Parmi celle-ci se trouvent les Nanoarchaeota, les *Nanohaloarchaea*, les ARMAN, et deux nouvelles lignées, nommées 'Aenigmarchaeota' et 'Diapherotrites'. Les auteurs proposent d'élever ces quatre dernières lignées au rang de phyla, les *Nanohaloarchaea* prendraient alors le nom de 'Nanohaloarchaeota' et les ARMAN deviendraient 'Parvarchaeota'. Cet ensemble de 5 phyla est nommé d'après leur acronyme, « DPANN ». D'un point de vue phylogénétique, cette étude présente plusieurs lacunes. Tout d'abord, l'arbre des archées présenté n'est pas raciné, il est donc difficile de conclure quant à la monophylie des deux groupes proposés « TACK » + Euryarchaeota et « DPANN ». Ensuite, de nombreuses études précédentes ont montré que les trois lignées de petites archées Nanoarchaeota, *Nanohaloarchaea* et ARMAN avaient une évolution rapide et de type réduction de génome et, par conséquent, étaient en cela particulièrement sujettes à l'attraction de longues branches (Brochier et al. 2005; Baker et al. 2010; Narasingarao et al. 2012). Elles doivent donc être traitées avec beaucoup de précautions, particulièrement dans des études à large échelle où les artefacts peuvent être nombreux. Enfin, l'élévation au rang de phylum de plusieurs lignées doit aussi être soumise à la

plus grande rigueur étant donné les deux précédents arguments (Simonetta Gribaldo and Brochier-Armanet 2012).

3. Conclusion

Les relations anciennes entre les différentes lignées d'archées ne sont donc toujours pas complètement résolues. L'utilisation, à peu de chose près, du même set de protéines comme marqueurs de la phylogénie quelle que soit la question posée (résolution de la phylogénie des archées, position de la racine, relations avec les autres domaines...) et l'absence d'analyses clairement orientées dans ce but sont peut être des débuts de réponses quant à la non résolution de ces points essentiels dans la phylogénie et la connaissance des archées, mais aussi de l'ensemble du vivant. C'est dans ce contexte que se place mon travail de thèse, par la recherche de nouveaux marqueurs de la phylogénie des archées et de la racine de l'arbre des archées par une approche phylogénomique.

Objectifs

Comme je l'ai montré dans l'Introduction, la compréhension des relations de parenté au sein des archées et entre les archées et les deux autres domaines du vivant, bactéries et eucaryotes, est essentielle pour comprendre l'histoire évolutive de l'ensemble du vivant. L'augmentation du nombre de génomes d'archées séquencés et, surtout, la diversité qu'ils recouvrent, ont ouvert la porte à des études massives sur ces questions évolutives.

Mon travail de thèse s'inscrit dans ce contexte, avec pour objectif une amélioration de la connaissance de la phylogénie des archées, et particulièrement en ce qui concerne les relations anciennes dans ce domaine. Ce travail s'est articulé autour de deux axes:

Objectif 1 : La recherche de nouveaux marqueurs pour l'inférence de la phylogénie des archées.

Peut-on trouver de nouvelles protéines pour résoudre la phylogénie des archées ? L'histoire évolutive des protéines informationnelles, utilisées jusqu'à maintenant comme marqueurs de la phylogénie des archées, est-elle représentative de l'histoire évolutive des organismes ou donne-t-elle un signal particulier ? Peut-on améliorer la résolution de la phylogénie des archées ? De nouveaux marqueurs peuvent-ils apporter des éléments de réponse concernant la position phylogénétique de certaines lignées, non résolue ou en débat, comme, par exemple, les relations au sein des méthanogènes classe I ; entre les méthanogènes classe II et les *Halobacteriales* ; la position des ordres de crenarchées proposés récemment, *Acidilobales* et *Fervidicoccales* ; la position des lignées de d'archées de taille nanométrique, *Nanoarchaeota*, *Nanohaloarchaea* et le groupe ARMAN ; ou les relations entre les différents phyla d'archées ; etc. ?

Pour cela, nous avons cherché de nouveaux marqueurs pour la phylogénie des archées par des approches de phylogénomique, puis inféré une phylogénie d'une part à partir des marqueurs sélectionnés et d'autre part en les associant avec les marqueurs de type informationnel utilisés dans les études précédentes (protéines ribosomiques, sous-unités de l'ARN polymérase et facteurs de transcription s). Cette étude fait l'objet du Chapitre 1.

Objectif 2 : La recherche de la racine de l'arbre des archées grâce à des homologues bactériens.

Parmi les marqueurs utilisés pour inférer la phylogénie des archées, certains d'entre eux ont-ils des homologues bactériens pouvant être utilisés pour raciner cette phylogénie (i.e. des homologues bactériens issus de l'ancêtre commun de tous les organismes et non pas de transferts horizontaux depuis des archées) ? Peut-on placer la racine (i.e. la première divergence au sein des archées) grâce à ces homologues bactériens avec un support suffisant ? Si oui, où se place cette

racine? Quelles relations entre les phyla d'archées observe t-on ? Quelles sont les implications sur la taxonomie des archées ?

Pour répondre à ces questions, nous avons recherché parmi les marqueurs utilisés pour inférer la phylogénie des archées (nouveaux et anciens) ceux qui avaient des homologues bactériens suffisamment conservés. Après une analyse fine des phylogénies individuelles contenant les homologues bactériens, nous avons sélectionné un sous-ensemble de marqueurs dédiés à cette étude, et inféré les phylogénies afin de raciner l'arbre des archées. Cette étude fait l'objet du Chapitre 2.

L'analyse des génomes d'archées nécessaires pour atteindre l'objectif 1 a produit une énorme masse d'information sur des protéines d'archées (leur répartition taxonomique, présence de transferts horizontaux, etc.). Cette information peut être exploitée pour comprendre l'évolution de ce domaine. Au cours de ma thèse, j'ai analysé l'histoire évolutive d'une de ces protéines, la protéine DnaJ-Fer, présente uniquement chez les thaumarchées et les plantes, et du système chaperonne DnaK-GrpE-DnaJ avec lequel la protéine DnaJ-Fer interagit. Ce travail est présenté dans le Chapitre 3.

Matériels et Méthodes

Les trois Chapitres contenus dans cette thèse sont le résultat d'une analyse globale des génomes de *Cenarchaeum symbiosum* et *Nitrosopumilus maritimus*, deux Thaumarchaeota, et de 'Candidatus Caldiarchaeum subterraneum', seul représentant du phylum Aigarchaeota. Le but premier de ces analyses était de mettre en évidence des protéines portant un signal pouvant servir à la résolution de la phylogénie des Archaea, point qui sera traité dans le Chapitre 1. Le choix de deux génomes de thaumarchées et de l'aigarchée nous a semblé le plus intéressant pour rechercher de nouveaux marqueurs pour la phylogénie des archées. D'une part, ils partagent des gènes avec les crenarchées et/ou les euryarchées (Brochier-Armanet et al. 2008; Nunoura et al. 2011), ce qui nous permettra de trouver de nouvelles protéines présentes dans les différents phyla. D'autre part, la position phylogénétique de ces phyla récemment proposés n'est pas résolue, et l'utilisation de marqueurs présents dans leurs génomes pourrait nous aider à répondre à cette question. Le Chapitre 2 traite de la position de la racine de l'arbre des Archaea, donc des relations entre phyla. Le Chapitre 3 présente une histoire évolutive originale, celle de la protéine DnaJ-Fer et du système chaperonne DnaK-GrpE-DnaJ, mettant en jeu plusieurs transferts horizontaux entre les trois domaines du vivant. L'analyse des génomes de *N. maritimus*, *C. symbiosum* et 'Ca. Caldiarchaeum subterraneum' a permis de souligner de nombreux autres cas de protéines avec des histoires évolutives potentiellement intéressantes, avec transferts horizontaux, duplications et/ou pertes de gènes, qui n'ont pu être traités en détails dans cette thèse.

J'ai d'abord étudié l'histoire évolutive de chaque protéine codée dans ces trois génomes par des analyses phylogénétiques préliminaires. J'ai ainsi pu sélectionner les protéines intéressantes pour le Chapitre 1, puis pour le Chapitre 2. Les méthodes utilisées pour l'analyse initiale et une partie des méthodes utilisées pour les Chapitres 1 et 2 sont donc communes ; je les expose ici.

A. Analyses phylogénétiques des protéines codées dans les génomes de *C. symbiosum*, *N. maritimus* et '*Ca. Caldiarchaeum subterraneum*'.

La première étape d'analyse a été la construction de la phylogénie de chacune des protéines codées dans les trois génomes de référence. Cette étape a été faite en collaboration avec Philippe Deschamps, qui avait mis en place précédemment une procédure qui génère des phylogénies de façon automatique pour les protéines choisies en utilisant les séquences similaires extraites d'une banque de données locale.

1. Construction d'une banque de données locale

La banque de données locale est composée des protéines codées par 92 génomes complets d'Archaea (correspondant à un représentant par espèce) disponibles en avril 2011; des protéines codées par les génomes complets de 297 bactéries et de 122 eucaryotes choisis parmi les différents phyla afin de représenter, dans la mesure du possible, la diversité actuellement connue et séquencée de ces deux domaines (Annexe 1, Supplementary Table S1). Ces séquences protéiques d'archées et de bactéries ainsi que la plupart de celles des eucaryotes ont été obtenues directement à partir de la banque de données GenBank (Benson et al. 2013), accessible sur le site du NCBI (<http://www.ncbi.nlm.nih.gov/>). Un certain nombre de génomes complets d'eucaryotes n'étant pas disponibles dans GenBank, ils ont pu être obtenus à partir d'autres banques, telles que le JGI (<http://genome.jgi-psf.org/>), et des banques spécialisées pour le génome de différents organismes (*Cyanidioschyzon merolae* : <http://merolae.biol.s.u-tokyo.ac.jp/> ou *Galdieria sulphuraria* : <http://genomics.msu.edu/galdieria/>). La liste de ces génomes est donnée dans l'Annexe 1, Supplementary Table S1.

2. Génération des phylogénies préliminaires

Le nombre de protéines étudiées et conservées à chaque étape est indiqué dans le Tableau 3. Les génomes de *N. maritimus* et de *C. symbiosum* étant assez proches en termes de similarité de séquences (e.g. 94% d'identité entre leurs ARN de la petite sous unité du ribosome (SSU rRNA)) et de contenu en gènes, nous avons décidé de ne conserver que la protéine issue du génome de *N. maritimus* dans le cas où une protéine homologue proche existait chez *C. symbiosum*. Pour cela, nous avons fait une recherche de similarité par BLASTp (Altschul et al. 1990) entre chaque paire de protéines homologues dans ces deux génomes, et si deux protéines présentaient un pourcentage d'identité supérieur ou égal à 60%, nous n'avons conservé que la protéine de *N. maritimus* (Tableau

3). Aux protéines ainsi triées, nous avons rajouté celles de '*Ca. Caldiarchaeum subterraneum*'.

A partir de chaque protéine des trois génomes étudiés, nous avons réalisé une recherche de similarité sur la banque de données locale grâce au logiciel BLASTp (Altschul et al. 1990). Nous avons gardé uniquement les 200 premiers résultats trouvés, ayant une E-value inférieure à 10^{-5} ; au dessus de cette valeur, nous considérons que, même si les séquences peuvent être homologues, elles sont trop distantes pour être utilisées dans notre étude. En effet, nous cherchons seulement les orthologues les plus proches de nos séquences de référence et si ceux-ci sont trop divergents, le signal évolutif ne sera pas assez fiable pour les utiliser comme marqueurs de la phylogénie des archaea. Suite à cette sélection, nous avons éliminé tous les jeux de données contenant moins de 4 séquences, puisqu'aucune information phylogénétique utile ne pourrait en être extraite (Tableau 3).

L'ensemble de ces jeux de données ont été alignés par le logiciel MAFFT (Katoh et al. 2005), avec les options par défaut. De précédents tests de différents logiciels d'alignements (ClustalW (Larkin et al. 2007), Muscle (Edgar 2004), Probcons (Do et al. 2005), T-Coffee (Notredame, Higgins, and Heringa 2000)), nous ont permis de comparer le rapport entre la qualité des résultats et le temps de traitement des données, et MAFFT nous semble être actuellement le programme le plus efficace dans la plupart des cas d'après ces deux critères. Certains cas complexes sont parfois mieux traités par d'autres logiciels (cf. Chapitre 3) mais dans une analyse massive de données, il était important d'utiliser la méthode qui donnait les meilleurs résultats de façon générale.

A partir de ces alignements, les positions dont l'homologie est relativement certaine ont été sélectionnées automatiquement grâce au logiciel BMGE (Criscuolo and Gribaldo 2010), avec les options par défaut. Ici encore, des tests effectués avec d'autres logiciels nous ont fait favoriser BMGE car il semblait être le plus efficace en termes de rapport qualité des résultats/temps de traitement.

Enfin, des phylogénies ont été inférées à partir de ces alignements grâce au logiciel FastTree version 2.1.4 (Price, Dehal, and Arkin 2010) implémentant une méthode de maximum de vraisemblance approximée, avec les options par défaut.

Génome	<i>N. maritimus</i>	<i>C. symbiosum</i>	' <i>Ca. Caldiarchaeum subterraneum</i> '
Nombres de jeux de données			
Total au départ (nombre de protéines codées dans le génome)	1796	2017	1704
Après BLAST réciproque entre <i>N. maritimus</i> et <i>C. symbiosum</i>	1796	674	1704
Après élimination des jeux de données contenant moins de 4 séquences	1381	626	1205
Total analysés phylogénétiquement	3212		

Tableau 3 : Nombre de protéines conservées à chaque étape de l'analyse des génomes complets de *N. maritimus*, *C. symbiosum* et '*Ca. Caldiarchaeum subterraneum*'.

3. Tri et sélection des phylogénies préliminaires

L'analyse des 3212 phylogénies inférées a été faite de façon manuelle avec l'aide du logiciel TreeReader développé dans notre équipe par Philippe Deschamps. Ce logiciel permet d'afficher des arbres phylogénétiques sur lesquels les noms d'espèces présents aux nœuds terminaux sont colorés automatiquement en fonction de leur taxonomie (d'après la classification utilisée par GenBank). L'autre avantage de cet outil est qu'il permet de charger un seul fichier contenant la totalité des arbres à étudier, au format Newick, et de passer d'un arbre à un autre d'un simple clic. Il est possible de raciner l'arbre observé sur le nœud choisi ou de faire permuter des branches. Cette facilité de lecture des arbres (à la fois en termes de manipulation et d'observation par la coloration automatique des noms d'espèces en fonction de leur taxonomie) a rendu plus efficace l'analyse de l'ensemble des phylogénies inférées par l'observation de chaque phylogénie.

Ainsi pour chaque arbre, j'ai noté dans un tableau ([Annexe 2](#)) une série d'observations. Tout d'abord, j'ai indiqué l'absence ou la présence de différents groupes taxonomiques, à commencer par les trois domaines Archaea, Bacteria et Eucarya, mais aussi à un niveau taxonomique plus détaillé, comme les différents phyla d'archées, ou les grands embranchements d'eucaryotes. Ensuite, j'ai inféré l'intérêt potentiel de la protéine comme marqueur de la phylogénie des Archaea. Pour cela, il fallait que la protéine soit présente dans un minimum de génomes d'Archaea (~10-15), que les séquences d'Archaea forment un groupe globalement monophylétique (afin de limiter les risques de transferts horizontaux et de paralogie ancienne) et que les groupes (particulièrement les ordres) dont la monophylie est décrite chez les Archaea le soient. J'ai aussi noté si de potentiels transferts

horizontaux étaient observables dans l'histoire évolutive de chaque protéine. Enfin, j'ai relevé un certain nombre d'observations, qui permettront une ré-analyse ciblée de certaines protéines pour aborder différentes problématiques. Les trois premiers points ont été notés de la façon la plus standardisée possible. Par exemple, l'absence d'un groupe est noté par un 0, sa présence par un 1 ; une protéine potentiellement marqueur de la phylogénie des Archaea est annotée « M ». Ce type d'annotation est précieux pour pouvoir faire différents tris automatiques sur ces données, et sortir la liste de protéines correspondante à une situation d'intérêt, même si elle n'avait pas été prévue au départ de l'analyse. Les analyses du Chapitre 2 en sont un parfait exemple, puisque ce projet n'avait pas été prévu lors du lancement de ce travail, mais a été possible grâce à un tri des familles protéiques sur le critère de la présence de bactéries dans les jeux de données.

B. Inférence de la phylogénie des Archaea : méthodes du Chapitre 1

1. Analyse des protéines d'intérêt pour l'étude de la phylogénie des Archaea

a. Sélection des protéines d'intérêt

A partir des informations collectées dans le tableau mentionné ci-dessus ([Annexe 2](#)) j'ai extrait une liste de protéines pouvant potentiellement être utilisées pour l'analyse de la phylogénie des Archaea. Comme précisé précédemment, il fallait que la protéine soit présente dans un minimum de génomes d'Archaea (~10-15), que les séquences d'archées soient globalement monophylétiques et que les groupes d'archées dont la monophylie est admise le soient. 209 protéines ont été sélectionnées à partir des génomes de *N. maritimus* et *C. symbiosum* (193 et 16 respectivement) et 55 à partir du génome de '*Ca. Caldiarchaeum subterraneum*' ([Tableau 4](#)) Ces chiffres excluent les protéines ribosomiques, les sous-unités de l'ARN polymérase et les facteurs de transcriptions y étant associés, bien connus et déjà utilisés dans de nombreuses études antérieures. Ces derniers ont fait l'objet d'une mise à jour au cours de mon travail à partir des données de Brochier-Armanet et collaborateurs en 2011 (Brochier-Armanet, Forterre, and Gribaldo 2011) et de Gribaldo et Brochier-Armanet en 2006 (Simonetta Gribaldo and Brochier-Armanet 2006), respectivement.

b. Construction d'une banque de données locale de génomes complets d'Archaea

L'objectif de cette analyse était de reconstruire la phylogénie globale des Archaea. De fait, nous avons décidé de ne travailler qu'à partir de séquences d'Archaea (un transfert horizontal vers un autre domaine ne pouvant modifier l'histoire évolutive au sein des Archaea) mais avec un échantillonnage le plus riche possible pour ce domaine. Pour cela, nous avons construit une banque de données locale nommée '*AllArchaea*', contenant toutes les protéines prédites de 129 génomes complets d'Archaea (avec un représentant par espèce) disponibles au début de cette analyse, en avril 2012 ([Annexe 1, Supplementary Table S3](#)). Les protéines de ces 129 génomes ont été collectées sur GenBank (Benson et al. 2012).

c. Construction des jeux de données

Pour chaque protéine, un nouveau jeu de données a été construit contenant les homologues identifiés chez les différentes espèces d'archées. Afin de couvrir au mieux la diversité des archées

lors de la recherche d'homologues par similarité avec BLASTp, je suis partie systématiquement de plusieurs séquences graines. Pour cela, j'ai sélectionné dans les phylogénies de départ (contenant des homologues bactériens et eucaryotes) une séquence de thaumarchée, d'euryarchée et de crenarchée dans le groupe monophylétique d'archées m'intéressant. La plupart du temps, la séquence de départ de *N. maritimus*, *C. symbiosum* ou '*Ca. Caldiarchaeum subterraneum*' a été sélectionnée, mais dans certains cas rares, j'ai préféré me concentrer sur un groupe de séquences paralogues à ma séquence d'origine. Les critères pour le choix d'une séquence graine étaient : une position phylogénétique fiable (pas de transfert horizontal de gènes (THG) ou de cas de paralogie flagrant, et pas de séquence isolée loin des autres archées) et une longueur de branche pas trop importante par rapport aux autres branches de l'arbre. Dans la mesure du possible, j'ai sélectionnée les séquences de *Methanobrevibacter ruminantium* et *Sulfolobus tokodaii* comme graine euryarchée et crenarchée respectivement. Si l'une de ces espèces était absente, j'ai sélectionné la séquence d'une espèce proche.

A partir de chacune de ces séquences graines, j'ai fait une recherche de similarité en local grâce au logiciel BLASTp (Altschul et al. 1997) sur la banque de données 'AllArchaea', avec les options par défaut, sauf pour le nombre de résultats et le nombre d'alignements visibles qui ont été réglés à 500. Nous avons considéré qu'au-delà de 500 résultats, soit les séquences n'étaient pas des homologues à notre protéine de départ, soit elles étaient trop éloignées pour la problématique de notre analyse, qui nécessite l'utilisation d'orthologues et doit exclure absolument les paralogues. Les résultats de BLAST ont été analysés manuellement afin de s'assurer de l'homologie des séquences retenues, et si possible de l'orthologie, pour la suite de l'analyse. Les séquences sélectionnées ont été extraites de la banque de données locale 'AllArchaea' grâce à un script de traitement de données codé par Philippe Deschamps, permettant de retrouver à partir d'un fichier de sortie BLAST au format xml les premiers résultats de BLAST jusqu'à une limite choisie chiffrée, et d'aller chercher les séquences correspondantes dans la banque de données utilisée pour la recherche de similarité. Un fichier contenant ces séquences au format fasta est alors créé.

Pour chaque protéine d'intérêt, j'ai donc, à cette étape, trois fichiers contenant les séquences issues des recherches de similarités depuis les séquences graines de Thaumarchaeota, d'Euryarchaeota et de Crenarchaeota. Ces trois fichiers ont alors été fusionnés en un seul, qui a été trié automatiquement pour éliminer les doublons et ne conserver qu'une seule copie de chaque séquence. A ce moment, nous avons un fichier fasta pour chaque protéine sélectionnée au départ comme marqueur potentiel.

d. Analyse phylogénétique des jeux de données

Alignement

Chaque jeu de données a été aligné grâce au logiciel MAFFT (Katoh et al. 2005) avec les options par défaut. Les fichiers d'alignements ont ensuite été convertis dans un format auquel est associé un fichier d'information contenant la fiche d'information GenBank de chaque séquence du fichier (fichier.ali et fichier.inf respectivement), ces deux fichiers étant lisibles par le logiciel MUST (Philippe 1993).

Phylogénies préliminaires

Chaque alignement a été nettoyé avec le logiciel BMGE (options par défaut) (Criscuolo and Gribaldo 2010) pour la suppression des positions non homologues, et une phylogénie a été inférée par maximum de vraisemblance grâce au logiciel PhyML version 3.0 (Guindon et al. 2010), avec la matrice de substitution LG (Le and Gascuel 2008) et une loi gamma à quatre catégories pour tenir compte de l'hétérogénéité de la vitesse d'évolution des sites (LG+G4). A cette étape, nous avons choisi d'utiliser le logiciel PhyML plutôt que le logiciel FastTree (approximation de maximum de vraisemblance) car il s'agit d'une méthode plus robuste, même si le temps de calcul est beaucoup plus long (il n'était pas envisageable de l'utiliser pour traiter plus de 3000 jeux de données de l'analyse préliminaire).

A partir des phylogénies produites par PhyML il a été possible de supprimer les redondances entre les jeux de données. En effet, il est possible par exemple que deux paralogues présents chez les Thaumarchaeota aient été traités séparément, alors que chez le reste des archées il n'existe qu'une seule copie de cette protéine. Dans ce cas, nous avons décidé de ne garder qu'une seule copie du jeu de données et qu'une seule copie de la protéine chez les Thaumarchaeota. Après cette étape, restaient donc 181 jeux de données pour les génomes de *N. maritimus* et *C. symbiosum* (173 et 8 respectivement) et 55 pour le génome de '*Ca. Caldiarchaeum subterraneum*'.

Répartition taxonomique des séquences

Enfin, à partir de chaque fichier fasta, j'ai reporté automatiquement le nombre de copies de la protéine présent dans chaque espèce d'Archaea afin d'avoir un indicateur du nombre de paralogues ou de copies en double dans le jeu de données. Ce travail a été fait grâce au script « TaxRep.pl », que j'ai codé pour compter le nombre de séquences d'une espèce (à partir de son numéro TaxID) dans un fichier fasta. Le résultat obtenu est un tableau dans lequel chaque ligne correspond à une espèce et chaque colonne à un fichier fasta. Le croisement entre une ligne et une

colonne donne le nombre de séquences de cette espèce dans le fichier ([Annexe 3](#)).

Analyse de l'alignement

Chaque jeu de données aligné a été analysé grâce au logiciel MUST (Philippe 1993) qui permet de faire des aller-retour entre l'alignement et des phylogénies inférées par Neighbour Joining (Saitou and Nei 1987). Cette étape a permis de supprimer les séquences très divergentes qui auraient pu rester, de sélectionner un groupe de paralogues dans les cas de duplication ancienne et de choisir entre deux copies paralogues pour une espèce ou pour un taxon si la duplication s'est faite après le dernier ancêtre commun des Archaea. Elle permet aussi de mettre en évidence d'éventuels transferts horizontaux de gènes entre Archaea, et donc d'écarter les séquences xénologues de l'analyse finale (dans ces cas particuliers, j'ai gardé une trace de l'alignement contenant les séquences issues de transferts horizontaux). Enfin, j'ai écarté un certain nombre de jeux de données dont la qualité n'était pas suffisante pour être utilisés pour l'analyse de la phylogénie des Archaea (irrésolution trop importante, doute sur des cas de THG ou de paralogie...) J'ai fait ce travail de nettoyage de l'alignement en tenant compte du tableau d'information de départ ([Annexe 2](#)), des phylogénies inférées par FastTree contenant les homologues bactériens et eucaryotes, des phylogénies PhyML ne contenant que les séquences d'archées, et du tableau de répartition taxonomique ([Annexe 3](#)). A la fin du nettoyage de l'alignement, il ne devait rester qu'une seule séquence par génome d'archée, donc au plus 129 séquences par jeu de données. Le nombre de jeux de données conservés est indiqué dans le [Tableau 4](#). La liste des protéines retenues est indiquée dans l'[Annexe 1, Supplementary Table S2](#).

Génome	<i>N. maritimus</i>	<i>C. symbiosum</i>	' <i>Ca. Caldiarchaeum subterraneum</i> '
Nombres de jeux de données			
Etudiés comme marqueurs potentiels	193	16	55
Après suppression des paralogues chez les Thaumarchaeota	173	8	55
Conservés comme marqueurs	153	4	43
Total des jeux de données retenus	200		

Tableau 4 : Nombre de protéines conservées à chaque étape de l'analyse des marqueurs potentiels depuis les génomes complets de *N. maritimus*, *C. symbiosum* et '*Ca. Caldiarchaeum subterraneum*'.

Les 200 jeux de données retenus ont été réalignés avec le logiciel MAFFT (Kato et al.

2005) avec options par défaut afin d'éliminer le bruit qui avait pu être créé dans l'alignement à cause des séquences très divergentes ou non homologues. Les positions homologues conservées utilisées pour l'inférence des phylogénies ont ensuite été choisies de manière semi-automatique pour une meilleure expertise grâce au programme NET du package MUST (Philippe 1993).

Phylogénies individuelles définitives

Enfin, pour chaque jeu de données, une phylogénie a été inférée par maximum de vraisemblance grâce au logiciel Treefinder version octobre 2008 (Jobb, von Haeseler, and Strimmer 2004) avec le modèle d'évolution LG+G4.

e. Analyse fonctionnelle des 200 nouveaux marqueurs

Pour avoir une idée du rôle de ces protéines dans la cellule, nous avons utilisé les catégories fonctionnelles de la banque COG (R L Tatusov, Koonin, and Lipman 1997; Roman L Tatusov et al. 2003). L'assignation à ces catégories fonctionnelles a été faite par une recherche de similarité par BLASTp contre la banque COG (<http://www.ncbi.nlm.nih.gov/COG/> (Roman L Tatusov et al. 2003)) à partir de chaque séquence protéique. Si le premier résultat de BLAST avait une E-value inférieure à 10^{-5} , j'ai assigné sa catégorie fonctionnelle COG à la séquence testée (Annexe 1, Supplementary Table S2).

2. Mise à jour des jeux de données de protéines informationnelles

Les analyses précédentes de la phylogénie des archées par Brochier-Armanet et collaborateurs ont été faites principalement sur les protéines de type informationnel telles que les protéines ribosomiques (Brochier-Armanet, Forterre, and Gribaldo 2011) ou les protéines formant les sous-unités de l'ARN polymérase (Simonetta Gribaldo and Brochier-Armanet 2006), comme expliqué dans l'Introduction. Dans cette nouvelle analyse de la phylogénie des Archaea, nous avons voulu réutiliser cet ensemble de protéines en plus des nouvelles données apportées par notre analyse des génomes complets de *N. maritimus*, *C. symbiosum* et '*Ca. Caldiarchaeum subterraneum*'. La liste des protéines issues de ces deux précédentes études est indiquée dans l'Annexe 4. Les jeux de données correspondants ont été mis à jour pour qu'ils présentent le même échantillonnage taxonomique que les 200 marqueurs que j'ai identifiés. De la même façon, les positions homologues conservées utilisées pour l'inférence phylogénétique ont été choisies manuellement avec le logiciel MUST.

3. Inférence de la phylogénie globale des Archaea

a. Construction des supermatrices

Le but de ce travail était d'une part d'améliorer la résolution de la phylogénie des Archaea, et d'autre part d'étudier le signal phylogénétique contenu dans d'autres protéines que les protéines informationnelles utilisées classiquement. Pour cela, nous avons utilisé une approche de supermatrice de caractères en concaténant différents alignements afin d'avoir un maximum de signal pour l'inférence phylogénétique. Dans cette optique plusieurs concaténations ont été construites : une première avec l'ensemble des 200 jeux de données issus de l'analyse présentée ici, une deuxième avec l'ensemble des protéines ribosomiques et une troisième contenant les sous-unités de l'ARN polymérase et les facteurs de transcription. Enfin, deux supermatrices ont été construites à partir de l'ensemble des marqueurs : une première contenant la totalité des 273 jeux de données et, afin de tester l'impact des données manquantes, nous avons aussi construit une supermatrice en limitant le nombre d'espèces manquantes autorisé à 10 par jeux de données, appelée « Matrice 10 ». Si plus de 10 espèces étaient absentes, le jeu de données n'était utilisé pour la construction de la supermatrice. Les différentes matrices construites et analysées sont indiquées dans le [Tableau 5](#).

	Nouveaux marqueurs	Protéines ribosomiques	ARN polymérase et facteurs de transcription	Matrice totale	Matrice 10
Nombre de positions	48 904	6228	2970	58 102	35 589
Nombre de jeux de données	200	57	16	273	179

Tableau 5 : Les différentes supermatrices construites pour les analyses phylogénétiques.

b. Désaturation

La saturation du signal phylogénétique dans certains marqueurs peut être responsable de l'irrésolution d'une phylogénie, ou du moins de l'irrésolution de la position d'espèces évoluant particulièrement vite (Philippe and Laurent 1998). Cette saturation est liée à des temps d'évolution très longs et/ou à des taux de mutation rapides, qui vont effacer la trace des mutations précédentes, donc du signal évolutif le plus ancien. Les protéines n'évoluant pas toutes à la même vitesse et, au sein de chaque protéine, chaque position pouvant être soumise à une pression de sélection

différente, une possibilité pour diminuer ce biais est de sélectionner les marqueurs et/ou les sites ayant une vitesse d'évolution plus lente. Les deux méthodes ont été employées dans cette étude afin d'essayer de résoudre certaines nœuds problématiques de la phylogénie des archées.

Désaturation par sélection de sites

La désaturation par sélection de sites a été faite par une méthode de « slow-fast » mise en place par Brinkmann et Philippe (Brinkmann and Philippe 1999). Le calcul des vitesses d'évolution par site et la construction des différentes matrices a été faite grâce au logiciel SlowFaster (Kostka et al. 2008). Cette méthode consiste à calculer pour chaque position de l'alignement la vitesse d'évolution du site (nombre de changements survenus par position) dans différents ensembles de séquences définies préalablement. Le taux d'évolution d'une position est estimé par la somme des changements survenus au sein des différents groupes de séquences sélectionnés. Les sites sont ensuite assemblés en différentes matrices : la première contenant les sites ayant le taux d'évolution le plus lent, auxquels sont ajoutés les sites ayant un taux d'évolution de plus en plus rapide. Ainsi, la matrice S_0 contient les sites invariants au sein de chaque ensemble de séquence, la matrice S_1 contient en plus les sites pour lesquels un seul changement a été mis en évidence sur l'ensemble des sous-groupes, et ainsi de suite.

Ici, nous avons choisi d'appliquer cette méthode sur la supermatrice contenant à la fois les nouveaux marqueurs et les protéines informationnelles contenant les jeux de données dans lesquels manquent au plus 10 espèces, ainsi que sur une sélection de 81 espèces (Annexe 1, Supplementary Table S3) correspondant à une liste réduite d'archées avec un choix équilibré de représentants de chaque groupe d'Archaea. 48 matrices ont été construites. Une deuxième série de matrices a été calculée à partir de la même supermatrice de départ mais sans les lignées d'archées nanométriques (*Nanoarchaeum equitans*, le groupe ARMAN (i.e. 'Candidatus Micrarchaeum acidiphilum' ARMAN-2, 'Candidatus Parvarchaeum acidiphilum' ARMAN-4, 'Candidatus Parvarchaeum acidophilus' ARMAN-5) et les Nanohaloarchaea (i.e. 'Candidatus Haloredivivus sp. G17', 'Candidatus Nanosalina sp. J07AB43' et 'Candidatus Nanosalinarum sp. J07AB56'). 47 matrices ont été construites. Pour les inférences phylogénétiques, nous avons choisi d'utiliser seulement les matrices séparées d'une taille d'au moins 500 positions dans un souci temps de calcul /gain de signal, ce qui représente 35 matrices. La liste des matrices utilisées et le nombre de positions qu'elles contiennent sont indiquées dans l'Annexe 5. Seuls les résultats sans les sept espèces placées sur des longues branches sont présentés dans le Chapitre 1.

Désaturation par sélection de gènes

Si la vitesse d'évolution est hétérogène entre les sites au sein d'une protéine, elle l'est également d'une protéine à l'autre. Nous avons voulu prendre en compte de cette hétérogénéité afin de désaturer la supermatrice totale et essayer de faire ressortir le signal porté par des protéines évoluant lentement. Pour cela, nous avons mis en place la méthodologie suivante : pour chaque protéine, un arbre individuel a été inféré par maximum de vraisemblance (avec Treefinder, cf. Matériels et Méthodes B.1.d), les longueurs de branches dans un arbre en maximum de vraisemblance étant le reflet de la vitesse d'évolution estimée de la protéine, nous avons utilisé ces longueurs de branches pour estimer une vitesse d'évolution moyenne de cette protéine. Ensuite, tout comme pour la désaturation par sélection de sites, nous avons construit différentes matrices, la première contenant les protéines évoluant le plus lentement, auxquelles sont ajoutés matrices après matrices, les protéines évoluant plus rapidement.

Concrètement, pour chaque protéine, j'ai extrait des matrices de distances entre chaque paire de séquences grâce au script Newick2Additive.jar développé par Alexis Criscuolo [Communication personnelle]. À partir de ces matrices de distances, j'ai calculé une « distance moyenne » pour l'arbre, en moyennant les distances séparant des couples d'espèces choisies préalablement. Pour cela, nous avons choisi de prendre en compte la distance entre deux groupes d'euryarchées (*Thermococcales* et *Methanomicrobiales*), entre deux groupes de crenarchées (*Sulfolobales* et *Thermoproteales*), entre chacun de ces groupes et les thaumarchées ou l'aigarchée, et entre chaque groupe d'euryarchées et chaque groupe de crenarchées. Pour chaque groupe, l'espèce présente dans le plus de jeux de données était utilisée comme référence, si elle était absente la suivante la plus représentée était utilisée. Au total, au maximum 10 distances ont donc été moyennées. Ce travail a été effectué automatiquement grâce à un script que j'ai mis en place « DistTree.pl », à partir des listes d'espèces utilisées et des matrices de distances extraites des arbres individuels. Enfin, nous avons construit, à partir des 273 jeux de données, 28 supermatrices, la première contenant les 10 jeux de données ayant les vitesses d'évolution les plus lentes, la seconde, les 20 jeux de données les plus lents et ainsi de suite ([Annexe 6](#)).

c. Inférences phylogénétiques

A partir de ces alignements, des phylogénies ont été inférées par différentes méthodes : en approximation de maximum de vraisemblance grâce au logiciel FastTree 2.1.4 (Price, Dehal, and Arkin 2010), en maximum de vraisemblance grâce au logiciel RaxML version 7.2.4 (Stamatakis 2006; Stamatakis, Hoover, and Rougemont 2008) avec des valeurs de soutien obtenues en Rapid

Bootstrapping et avec PhyML version 3.1 (Guindon et al. 2010), et par une méthode bayésienne avec Phylobayes version 3.3 (Lartillot, Lepage, and Blanquart 2009). Le calcul des phylogénies à partir de supermatrices aussi importantes peut prendre énormément de temps, et avec des différences très importantes selon les méthodes. Certains calculs sont donc toujours en cours, et seule une partie des résultats sont présentés dans le Chapitre 1.

Nous avons aussi fait différentes sélections d'espèces sur ces supermatrices afin d'essayer de résoudre des questions particulières de la phylogénie des Archaea. D'une part en enlevant les lignées correspondants aux archées nanométriques qui se positionnent sur des longues branches, qui peuvent être soumises à des artefacts tels que l'attraction de longues branches (LBA) (Philippe and Laurent 1998), et en les réinsérant une par une. D'autre part, en ne sélectionnant que les euryarchées d'une part et les crenarchées d'autre part permet d'obtenir un meilleur signal au sein de chaque phylum.

Le détail des méthodes et des paramètres utilisés est expliqué dans l'article du Chapitre 1.

C. Racinement de l'arbre des Archaea : méthodes du Chapitre 2

1. Génération des jeux de données pour le racinement de l'arbre des Archaea

Certaines des protéines que nous avons utilisées pour l'inférence de la phylogénie des Archaea ont des homologues bactériens. Les protéines ribosomiques, bien sûr, mais aussi certains des 200 nouveaux marqueurs mis en évidence dans cette étude. Nous avons donc utilisé ces homologues bactériens comme groupe extérieur pour essayer de raciner la phylogénie des Archaea.

a. Sélection des protéines d'intérêt

La présence d'homologues bactériens avait déjà été notée lors de l'analyse globale de l'ensemble des phylogénies des trois génomes de *N. maritimus*, *C. symbiosum* et '*Ca. Caldiarchaeum subterraneum*'. Il a suffi de sélectionner parmi les 200 marqueurs de la phylogénie des Archaea, les protéines ayant des homologues bactériens. Il y en avait 81.

b. Construction d'une banque de données locale de génomes complets de bactéries

L'objectif de cette analyse était de raciner la phylogénie des Archaea avec les homologues bactériens pour les marqueurs déjà analysés. Pour cela, nous avons construit une banque de données locale nommée 'BacteriaRacine', contenant les protéines prédites dans 117 génomes complets de bactéries, correspondant à environ cinq représentants par phylum bactérien. Nous avons choisi les représentants de chaque phylum en prenant les cinq génomes contenant le plus de protéines (un représentant par espèce et par genre, dans la mesure du possible) parmi les génomes complets disponibles au début de cette analyse en février 2013 ([Annexe 7, Supplementary Table S1](#)).

c. Construction des jeux de données pour l'analyse des marqueurs potentiels

Pour chacun des 81 marqueurs sélectionnés, nous avons fait une recherche de similarité par BLASTp (Altschul et al. 1997) sur la banque 'BacteriaRacine' à partir de la protéine de départ (en utilisant les homologues de *N. maritimus*, *C. symbiosum* ou '*Ca. Caldiarchaeum subterraneum*' comme graine dans la mesure du possible). Les recherches de similarité et l'analyse des résultats de BLAST ont été faites de la même façon que l'analyse des résultats de BLAST sur la banque 'AllArchaea' dans l'étape précédente de cette étude (cf. Matériels et Méthodes B.1). En cas de doute sur le fait que certains homologues bactériens n'aient pas été détectés par BLAST à partir de

la séquence d'archée, j'ai refait une recherche de similarité à partir de la première séquence bactérienne trouvée. Une fois que la sélection des séquences à garder pour l'analyse a été faite, la génération des jeux de données finaux a été réalisée en appliquant le même protocole que dans la première analyse (cf. Matériels et Méthodes B.1.c).

d. Analyse des jeux de données

Les jeux de données ont été analysés de la même façon que ceux lors de la première analyse, expliquée dans la partie Matériels et Méthodes B.1.d. A la fin du nettoyage de l'alignement, il ne devait rester qu'une seule séquence par espèce d'archée ou de bactérie, donc au plus 246 séquences par jeu de données. Le nombre de jeux de données conservés est indiqué dans le [Tableau 6](#). La liste des 38 protéines conservées est indiquée dans l'[Annexe 7, Supplementary Table S2](#) et dans l'[Annexe 8](#). La répartition taxonomique de ces 38 protéines parmi les 246 génomes espèces est donnée dans l'[Annexe 9](#).

Génome	<i>N. maritimus</i>	<i>C. symbiosum</i>	' <i>Ca. Caldiarchaeum subterraneum</i> '
Nombre de jeux de données			
Protéines ayant des homologues bactériens.	58	1	22
Jeux de données conservés pour le racinement de la phylogénie des Archaea	35	1	9
Total des jeux de données conservés pour l'analyse	38		

Tableau 6 : Nombre de jeux de données conservés à chaque étape de l'analyse visant à raciner l'arbre des Archaea avec les homologues bactériens.

2. Mise à jour des jeux de données de protéines informationnelles pour le racinement de l'arbre des Archaea

De la même façon que les jeux de données de protéines informationnelles ont été utilisés après mise à jour pour l'analyse de la phylogénie globale des Archaea, nous avons complété ces mêmes jeux de données en ajoutant les homologues bactériens des espèces que nous avons sélectionnées. Par contre, seules les protéines ribosomiques ont été utilisées ici, car peu de sous-unités de l'ARN polymérase archéenne ont des homologues bactériens, et pour les quelques sous-unités pour lesquelles c'est le cas, ces homologues sont très divergents, donc très peu de signal serait réellement exploitable. 32 protéines ribosomiques d'archées ont des homologues bactériens

bien conservés et dont le signal peut être utilisé. La liste de ces protéines utilisées est indiquée dans l'[Annexe 7, Supplementary Table S2](#).

3. Inférence phylogénétique et racinement de l'arbre des Archaea

a. Construction des supermatrices

De la même façon que les jeux de données utilisés pour la phylogénie globale des Archaea ont été concaténés en différentes supermatrices avec différents échantillonnages de jeux de données ou d'espèces, les jeux de données utilisés pour placer la racine de l'arbre des archaea l'ont été. La première supermatrice contient les 38 marqueurs issus de l'analyse précédente, une deuxième les 32 protéines ribosomiques, et une dernière contient l'ensemble de ces 70 jeux de données. Le but de cette analyse étant de placer la racine de l'arbre des archaea et non pas de résoudre la phylogénie intragroupe chez les archaea ou chez les bactéries, ces supermatrices ont été construites avec des listes de 81 archées et de 27 bactéries (108 espèces au totale) ([Annexe 7, Supplementary Table S1](#)) afin de limiter le temps de calcul, déjà très important pour des supermatrices de cette taille. Les différentes matrices construites et analysées sont indiquées dans le [Tableau 7](#). À partir de ces supermatrices, d'autres sélections d'espèces ont été choisies pour analyser l'effet des espèces placées sur des longues branches par exemple, tout comme dans l'analyse précédente.

	Nouveaux marqueurs	Protéines ribosomiques	Concaténation totale
Nombre de positions	6890	2560	9450
Nombre de jeux de données	38	32	70

Tableau 7 : Différentes concaténations utilisées pour les analyses phylogénétiques de racinement de l'arbre des archées.

b. Désaturation

Le positionnement de la racine de l'arbre des archaea nécessitant d'extraire un signal encore plus ancien que celui pour résoudre les relations intra-archées, les analyses de désaturation étaient aussi indiquées dans ce cas là. La méthodologie appliquée est la même que dans la première analyse.

Désaturation par sélection de sites

La désaturation par sélection de sites pour les jeux de données construits pour la résolution du placement de la racine de l'arbre des archées a été faite de la même façon que dans l'analyse précédente. 10 matrices ont été construites, chacune ayant un gain d'environ 1000 positions en taille par rapport à la précédente. La liste des matrices utilisées et le nombre de positions qu'elles contiennent sont indiquées dans l'[Annexe 10](#).

Désaturation par sélection de gènes

La désaturation par sélection de gènes a été faite de la même façon que dans l'analyse précédente, mais des distances entre trois espèces de bactéries et le reste des espèces d'archées déjà sélectionnées ont été utilisées en plus. Nous avons ainsi construit 10 matrices à partir des 70 jeux de données, la première contenant les 7 jeux de données ayant les vitesses d'évolution les plus lentes puis en ajoutant les jeux de données 7 par 7. ([Annexe 11](#))

c. Inférences phylogénétiques

De même que dans l'analyse (cf. Matériels et Méthodes B.3), des phylogénies ont été inférées avec différentes méthodes, particulièrement avec les logiciels FastTree (Price, Dehal, and Arkin 2010) et RAxML (Stamatakis 2006) dans les mêmes versions que précédemment, et MrBayes v3.2.1 (Ronquist et al. 2012) pour les analyses en inférence bayésienne.

Le détail des méthodes et des paramètres utilisés est expliqué dans l'article du Chapitre 2.

Chapitre 1 : La Phylogénie des archées, au-delà des protéines informationnelles.

A. Introduction

Résoudre la phylogénie des archées est, comme je l'ai expliqué dans l'Introduction, important pour la compréhension de l'histoire évolutive de ce domaine mais aussi de l'ensemble du vivant étant données les interactions entre les archées, leurs milieux et les autres êtres vivants. Les protéines ribosomiques (et par extension les protéines informationnelles, avec les sous-unités de l'ARN polymérase) sont utilisées de plus en plus couramment depuis le séquençage des premiers génomes d'archées pour établir cette phylogénie. Récemment, de grandes avancées dans les techniques de séquençage ont permis l'augmentation exponentielle de la quantité de génomes disponibles et, notamment, pour des taxons couvrant une part de plus en plus grande de la diversité des archées. De nombreuses études ont été conduites récemment (cf. Introduction) pour essayer d'inférer la phylogénie des archées à partir de ces données, mais les marqueurs utilisés restent toujours, pour une grande partie, des protéines informationnelles. De ces analyses ressort le fait que malgré leur qualité en tant que marqueurs phylogénétiques, ces protéines ne permettent pas de résoudre tous les nœuds de l'arbre des archées. De plus, les protéines ribosomiques, les sous-unités de l'ARN polymérase ou les ARNr interviennent tous dans des processus interconnectés dans la cellule et sont peut être porteurs d'un même signal phylogénétique, représentatif non pas de la phylogénie des organismes mais de ces systèmes en particulier.

Afin de répondre à ces questions, nous avons cherché de nouveaux marqueurs pour la phylogénie des archées. Pour cela, nous avons inféré et analysé la phylogénie de chaque protéine codée dans trois génomes (deux thaumarchées, *N. maritimus* et *C. symbiosum* et de l'aigarchée '*Ca. Caldiarchaeum subterraneum*'). Des jeux de données ont été construits en collectant les 200 premières séquences trouvées par BLAST contre une banque locale contenant tous les génomes d'archées disponibles, ainsi que des génomes de bactéries et d'eucaryotes représentatifs de la diversité de ces domaines. Un certain nombre de protéines a été écarté très rapidement, car ne celles-ci ne présentaient pas assez d'homologues. L'analyse des 3212 jeux de données restants est passée par plusieurs étapes de tri, jusqu'à la sélection de 200 protéines porteuses d'un signal suffisamment fiable (sans transfert horizontal de gène ou paralogie). Ces 200 protéines interviennent dans de nombreux processus cellulaires, pas uniquement informationnels. En effet, l'analyse des catégories fonctionnelles COG auxquelles elles sont assignées montre que 53 interviennent dans des voies métaboliques, 24 dans des processus cellulaires et de signalisation, et

24 appartiennent à des familles encore peu caractérisées. Une concaténation de ces 200 nouveaux marqueurs a été construite pour 129 génomes d'archées. Bien que 149 génomes soient disponibles, nous avons décidé de ne conserver qu'un seul génome par espèce dans la mesure où notre travail ne porte pas sur la phylogénie à un niveau taxonomique aussi faible. Parmi ces 129 espèces, 34 n'étaient pas représentées dans la dernière analyse de la phylogénie des archées (Brochier-Armanet, Forterre, and Gribaldo 2011), dont sept nouveaux genres. La phylogénie inférée à partir de cette première supermatrice montre un signal très proche de celui des protéines informationnelles utilisées couramment. Nous avons aussi actualisé ces jeux de données (57 protéines ribosomiques (Brochier-Armanet, Forterre, and Gribaldo 2011) et 16 sous-unités de l'ARN polymérase et facteurs de transcription associés (Simonetta Gribaldo and Brochier-Armanet 2006)), afin d'avoir le même échantillonnage taxonomique que pour les nouveaux marqueurs. Deux supermatrices ont été construites, l'une avec les protéines ribosomiques, l'autre avec le système de transcription. Le signal phylogénétique de ces trois supermatrices étant congruent, nous avons pu construire deux nouvelles supermatrices à partir de l'ensemble de ces 273 protéines. Une première contenait 179 jeux de données pour lesquels moins de 10 espèces sont manquantes (afin d'estimer l'impact des données manquantes) et la seconde contenait l'ensemble des jeux de données, avec un total de 58102 positions. Des méthodes de désaturation site par site et gène par gène ont été appliquées à ces supermatrices afin de réduire la saturation mutationnelle des jeux de données et pour essayer de replacer des espèces évoluant rapidement et sujettes aux artefacts d'attraction de longues branches. Pour la même raison, différents échantillonnages d'espèces ont été réalisés sur les archées de taille nanométrique (*N. equitans*, '*Ca. Micrarchaeum*', '*Ca. Parvarchaeum*' et *Nanohaloarchaea*).

Ce travail est présenté dans l'article « *Extending the conserved phylogenetic core of Archaea disentangles the evolution of the third domain of Life.* » (Petitjean, Deschamps, López-García, Moreira, Brochier-Armanet, en préparation).

B. Manuscrit de l'article 1 : «Extending the conserved phylogenetic core of Archaea disentangles the evolution of the third domain of Life. »

1 **Extending the conserved phylogenetic core of Archaea disentangles the evolution of the**
2 **third domain of Life.**

3
4 Céline Petitjean¹, Philippe Deschamps¹, Purificación López-García¹, David Moreira^{1,*}, Céline
5 Brochier-Armanet^{2,*}.

6
7 ¹UMR CNRS 8079, Unité d'Ecologie, Systématique et Evolution Université Paris-Sud, 91405
8 Orsay, Cedex, France.

9 ²Université de Lyon, Université Lyon 1, CNRS, UMR5558, Laboratoire de Biométrie et Biologie
10 Evolutive, 43 boulevard du 11 novembre 1918, F-69622 Villeurbanne, France. Tel.: +33 4 26 23
11 44 76; fax: +33 4 72 43 13 88.

12
13 Corresponding authors:

14 David Moreira (david.moreira@u-psud.fr) and Céline Brochier-Armanet (celine.brochier-
15 armanet@univ-lyon1.fr).

16
17

18
19

20 **Abstract.**

21 Seminal works aiming at studying the phylogeny of Archaea relied mainly on the analysis of the
22 RNA component of the small subunit of the ribosome (SSU rRNA). The resulting phylogenies
23 have provided interesting but partial information on the evolutionary history of the third domain of
24 life because SSU rRNA sequences do not contain enough phylogenetic signal. Therefore many
25 relationships, and especially the most ancient, remained elusive. Moreover, SSU rRNA
26 phylogenies can be heavily biased by tree reconstruction artifacts. The sequencing of complete
27 genomes allowed using protein markers as alternative to SSU rRNA and the ribosomal proteins
28 are now used routinely to study ancient phylogenies. Taking the opportunity of the recent burst
29 of archaeal complete genome sequences, we have carried out an in-depth phylogenomic
30 analysis. We have identified 200 new protein families that form a conserved phylogenetic core of
31 genes together with the ribosomal proteins and the subunits of the RNA polymerase. The
32 accurate analysis of these markers sheds new light on the evolutionary history of this domain.
33 We resolved a number of important relationships such as those among methanogens Class I.
34 Furthermore the use of desaturation approaches revealed that several relationships recovered in
35 recent analyses are the consequence of tree reconstruction artifacts and allowed replacing the
36 three very fast evolving lineages of nanosized archaea.

37

38 **Keywords.** Phylogenomics, *Methanopyrus kandleri*, Methanohoma, ARMAN, Nanoarchaeota,
39 Nanohaloarchaea, Horizontal gene transfer, mutational saturation, Slow-Fast method.

40

41 **Running title.**

42 The identification of 200 new conserved phylogenetic markers brings-up the phylogeny of
43 Archaea.

44 **Data deposition:** All sequence alignments used in this work are available upon request to the
45 corresponding authors.

46

47

49 **Introduction.**

50 The seminal work of Carl Woese and George Fox at the end of the 70's (Woese, Fox 1977) has
51 contributed to establish the RNA component of the small subunit of the ribosome (SSU rRNA) as
52 the gold standard to study the evolutionary relationships among living beings (and especially
53 among microorganisms), and indeed this marker was subsequently proven to be a powerful tool
54 for modern systematics and the exploration of microbial diversity. Among the most important
55 discoveries relying on the analysis of SSU rRNA sequences was the awareness that the living
56 world was divided into three domains (i.e. Archaea, Eucarya and Bacteria) (Woese, Fox 1977)
57 and that most of the biological diversity was represented by uncultured microorganisms (for a
58 recent review on the topic see (Lopez-Garcia, Moreira 2008)).

59 In the 90's, however, the question was asked whether the phylogenies based on SSU rRNA
60 sequences actually reflect the evolutionary history of organisms or, in other words, whether the
61 SSU rRNA is suitable to trace back the wealth of speciation events that have affected the
62 cellular lineages, especially the most ancient ones (Stiller, Hall 1999; Philippe, Germot 2000). In
63 fact, the phylogenetic signal carried by this molecular marker is too weak to resolve the deepest
64 nodes of the archaeal phylogeny, leading to largely unresolved trees (Robertson et al. 2005;
65 Cavicchioli 2011) but this is specific to neither Archaea nor SSU rRNA given that similar
66 situations have been reported for Bacteria and Eucarya and for other molecular markers (Roger
67 1999; Philippe et al. 2000; Brochier, Philippe 2002). The lack of phylogenetic signal can result
68 either from radiation, mutational saturation or from a combination of both (Gribaldo, Brochier
69 2009). Radiation is encountered when the diversification of the lineages under study occurred
70 too rapidly to be recorded at the molecular level, meaning that too few substitutions were fixed
71 between cladogenesis events. Conversely, mutational saturation results from the progressive
72 erasure of the most ancient phylogenetic signal by the accumulation of more recent substitutions
73 occurring at the same sites. As a consequence, in both cases the order of the speciation events
74 is hardly traceable by the phylogenetic analysis of present-day homologues of the studied
75 molecular marker. In addition, phylogenies based on SSU rRNA can be heavily affected by
76 several tree reconstruction artifacts, such as the Long Branch Attraction (LBA), which is due to
77 the heterogeneity of evolutionary rates among the studied sequences and leads to the grouping
78 of the fastest- and the slowest-evolving sequences in different parts of the tree (Felsenstein
79 1978). This artifact has been particularly well documented in the case of Eucarya and Metazoa
80 (see (Delsuc, Brinkmann, Philippe 2005) and references therein). Additional biases such as
81 those linked to compositional heterogeneity of sequences can also affect phylogenies (Delsuc,
82 Brinkmann, Philippe 2005). Prokaryotic SSU rRNA phylogenies are particularly sensitive to this
83 bias because the base composition of structural RNAs (e.g. SSU and LSU rRNA, tRNA, etc.) is
84 strongly correlated with the optimal growth temperature of the organisms (Woese et al. 1991;
85 Galtier, Lobry 1997). Finally, cases of horizontal gene transfer affecting SSU rRNA genes have
86 been reported (Yap, Zhang, Wang 1999; Bodilis et al. 2012; Kitahara, Miyazaki 2013).

87 Disentangle the deepest nodes of the Tree of Life, namely deciphering the relationships among
88 the main lineages within each of the three domains, is however crucial because it provides the
89 evolutionary frame indispensable to understand how the present-day diversity arose and how

90 biological features (e.g. metabolic processes, cellular structures, genomes, etc.) evolved all
91 along the diversification of Life (Gribaldo, Brochier 2009). To overcome the limited phylogenetic
92 signal carried by molecular markers, alternative approaches have been proposed and
93 successfully applied (Delsuc, Brinkmann, Philippe 2005). These included the development of
94 accurate evolutionary models overcoming some of the simplifying assumptions of the Markovian
95 models currently used in molecular phylogenetics (Lartillot, Philippe 2004; Le, Lartillot, Gascuel
96 2008; Groussin, Boussau, Gouy 2013) and reducing the risk of tree reconstruction artifacts
97 (Lartillot, Brinkmann, Philippe 2007). In parallel, the increase of computational power allowed
98 generalizing the use of statistical methods for phylogenetic inference (e.g. Maximum Likelihood
99 and Bayesian inference) that are less prone to tree reconstruction artifacts such as the LBA
100 (Delsuc, Brinkmann, Philippe 2005).

101 Beside methodological aspects, the past five years have witnessed a burst of large scale
102 genome sequencing projects covering an ever-growing part of the taxonomic diversity (including
103 the uncultured one) within the three Domains (Wu et al. 2009; Rinke et al. 2013). This windfall of
104 data provides valuable material to tackle complex evolutionary questions, allowing for instance
105 the selection of accurate taxonomic samplings targeting the slowly-evolving sequences within
106 each taxonomic group. Focusing on these sequences which are less susceptible to have
107 undergone multiple substitutions can help to reduce the mutational saturation level of the
108 datasets and thus limit the LBA (Delsuc, Brinkmann, Philippe 2005; Rodriguez-Ezpeleta et al.
109 2007). Last but not least, the availability of complete genomes has revolutionized phylogenetics,
110 shifting progressively to phylogenomics and thus from single-gene analysis towards the analysis
111 of hundreds of markers either through super-matrix or super-tree approaches (Delsuc,
112 Brinkmann, Philippe 2005). This allows combining the weak phylogenetic signal carried by each
113 individual molecular marker towards a stronger signal, and reducing the global level of noise
114 contained in the data by diluting the noise carried by each individual marker, providing that the
115 biases inherent to each marker are different. In return, phylogenomic approaches require a
116 crucial preliminary and time-consuming step aiming at identifying and selecting the orthologous
117 sequences of each studied marker.

118 Altogether these approaches have significantly improved our knowledge of ancient evolution. In
119 the case of Archaea, we have shown that components of various biological systems, in particular
120 transcription and translation, formed a conserved phylogenetic core that can be used to trace
121 back the evolutionary history of this domain (Brochier, Forterre, Gribaldo 2005; Gribaldo,
122 Brochier-Armanet 2006). From these analyses, a global picture of the evolutionary history of
123 Archaea is emerging (see (Brochier-Armanet, Forterre, Gribaldo 2011) and references therein).
124 These markers support the division of Archaea into four main phyla: Euryarchaeota,
125 Crenarchaeota, Thaumarchaeota (including the candidate phylum 'Aigarchaeota') and
126 Korarchaeota and has confirmed some relationships based on SSU rRNA analyses, among
127 which the grouping of Desulfurococcales and Sulfolobales within Crenarchaeota or the divide of
128 Euryarchaeota in three parts: a basal part containing Thermococcales and methanogens Class I
129 (i.e. Methanopyrales, Methanobacteriales and Methanococcales), an intermediate region
130 containing two groups: Thermoplasmatales and relatives (e.g. DHEV2, the uncultured marine
131 group II, etc.), and Archaeoglobales, and an apical region gathering Halobacteriales and
132 methanogens Class II (i.e. Methanocellales, Methanomicrobiales, Methanosarcinales, and

133 relatives). Despite these significant advances, the phylogeny of Archaea is far from being fully
134 resolved and a number of important nodes require further investigations (Brochier-Armanet,
135 Forterre, Gribaldo 2011). For instance, several lineages of uncultured nanosized archaea have
136 been discovered in very hot (i.e. Nanoarchaeota (Huber et al. 2002), highly acidic (i.e. ARMAN,
137 for Archaeal Richmond Mine Acidophilic Nano-organisms) (Baker et al. 2006)) and hypersaline
138 (i.e. Nanohaloarchaea (Narasingarao et al. 2012)) environments. Among them,
139 Nanohaloarchaea are widespread and could represent a significant part of the biodiversity of
140 hypersaline ecosystems (Narasingarao et al. 2012). Genome sequencing has revealed that
141 these lineages of tiny archaea are fast evolving and, accordingly, their phylogenetic position is
142 highly debated. Nanoarchaeota were first proposed to represent the earliest branching lineage
143 within Archaea, and therefore to represent a new phylum (Huber et al. 2002; Waters et al. 2003),
144 whereas later analyses suggested that Nanoarchaeota represent in fact a fast evolving
145 euryarchaeal lineage likely related to Thermococcales and that the basal branching observed in
146 some phylogenetic analyses was the consequence of a LBA (Brochier et al. 2005). However,
147 this point is still highly debated (see (Podar et al. 2013) and the comments of the reviewers
148 therein). Regarding ARMAN, a relationship with Thermoplasmatales was initially proposed
149 (Baker et al. 2006), but later analyses suggested that ARMAN could encompass distinct
150 lineages, among which one could be related to Nanoarchaeota (Brochier-Armanet, Forterre,
151 Gribaldo 2011). Finally, Nanohaloarchaea could represent a distant lineage related to
152 Halobacteriales (Narasingarao et al. 2012). However, this should be confirmed by additional
153 investigations. In addition to nanosized archaea, several other newly proposed lineages are
154 debated, such as the Acidilobales within Crenarchaeota (Prokofeva et al. 2009; Mardanov et al.
155 2010) or the candidate phylum 'Aigarchaeota' (Nunoura et al. 2011), (see (Brochier-Armanet,
156 Forterre, Gribaldo 2011) and references therein). Finally a number of important questions remain
157 to be resolved, among which the relationships among the four archaeal phyla and the root of
158 Archaea, the monophyly of some orders such as the Desulfurococcales, or the relationships
159 among methanogens Class I, and among methanogen Class II and Halobacteriales.

160 Using exhaustive comparative genomics and phylogenomics approaches, we have identified
161 200 new conserved proteins widely distributed among the archaeal phyla. Most of them are
162 involved in cellular activities other than transcription and translation. In combination with the
163 classical markers previously used to study the phylogeny of Archaea (namely, 73 transcription
164 and translation proteins), we have carried out phylogenetic analyses with different selections of
165 markers and taxa. Our results confirm many previous observations and clarify several uncertain
166 parts of the archaeal phylogeny, such as the monophyly of the methanogens Class I, the position
167 of Halobacteriales within the methanogens Class II or the refutation of the order Acidilobales.

168

169 **Materials and methods.**

170 *Dataset assembly.*

171 The proteomes of *Cenarchaeum symbiosum* (Hallam et al. 2006) and *Nitrosopumilus maritimus*
172 (Walker et al. 2010), two species of Thaumarchaeota, and of '*Candidatus* Caldiarchaeum
173 subterraneum', a composite genome assembled from a metagenomic library proposed to

174 represent a new candidate phylum tentatively called 'Aigarchaeota' (Nunoura et al. 2011), were
175 used as starting point to identify protein families of potential interest to study the phylogeny of
176 Archaea.

177 The 2017 proteins of *C. symbiosum* were compared to the 1796 proteins of *N. maritimus* using a
178 best-reciprocal blast-hit approach. Pairs of proteins displaying more than 60% of identity were
179 considered as members of the same family. This comparison led to the identification of 2470
180 protein families. Together with the 1704 proteins of 'Ca. Caldiarchaeum subterraneum', these
181 represented 4174 protein families containing at least one representative of Thaumarchaeota or
182 of the candidate phylum 'Aigarchaeota'.

183 Then, these protein families were used to query a local sequence database with blastp (Altschul
184 et al. 1997). This local database contained 92 archaeal, 297 bacterial and 122 eukaryotic
185 proteomes publicly available at the NCBI (<http://www.ncbi.nlm.nih.gov/>) (Supplementary Table
186 S1) that covered the taxonomic diversity of each domain of Life. The 200 first high-scoring
187 segment pairs with e-values lower than 10^{-5} were retrieved and added to the corresponding
188 protein family. At this step, 962 protein families containing less than four sequences were
189 discarded. The remaining 3212 datasets were aligned with MAFFT version 7 (default
190 parameters) (Kato, Standley 2013). The resulting alignments were trimmed using BMGE
191 (default parameters) (Criscuolo, Gribaldo 2010). Preliminary phylogenetic trees were inferred
192 with FastTree version 2 (JTT+CAT model) (Price, Dehal, Arkin 2010). The resulting 3212 trees
193 were visually inspected.

194 Among the 3212 datasets and beside ribosomal proteins and RNA polymerase subunits and
195 transcription factors, 236 corresponded to distinct protein families which could represent good
196 candidates to study the global phylogeny of Archaea because they produced phylogenetic trees
197 supporting strongly the monophyly of this domain and most of its orders. Then these protein
198 families were used to query a local database composed of all available archaeal complete
199 genome sequences (representing 129 different species, see Supplementary Table S3) with
200 blastp and at least one thaumarchaeotal, one crenarchaeotal and one euryarchaeotal sequence
201 as seed. All archaeal homologous sequences within each protein family were aligned using
202 MAFFT and the alignments were trimmed with BMGE. Maximum likelihood (ML) phylogenies
203 were inferred upon the trimmed alignments with PhyML version 3.1 (with the accurate NNI+SPR
204 and tlr options) (Guindon et al. 2010), with the Le and Gascuel (LG) model (Le, Gascuel 2008)
205 and a gamma distribution (Γ) to model the heterogeneity of evolutionary rates across sites (four
206 site categories). Branch robustness was estimated with the non-parametric bootstrap procedure
207 implemented in PhyML (100 replicates of the original alignments). The resulting trees were
208 visually inspected. 36 protein families of the original 236 were discarded at this step because
209 they present complex patterns of gene duplications and losses, and possible cases of horizontal
210 gene transfer among archaeal species that were not apparent when the taxonomic sampling
211 was restricted to 91 archaea. As a result, 200 conserved proteins were kept for final analyses
212 (Supplementary Table S2). When in-paralogous sequences were present in these 200 protein
213 families, we kept the slowest-evolving copy for the phylogenetic analyses.

214 In parallel, the 57 ribosomal proteins and 16 protein families corresponding to the 14 subunits of
215 the RNA polymerase and transcription factors used in previous studies on the phylogeny of

216 Archaea (Brochier, Forterre, Gribaldo 2004; Brochier-Armanet, Forterre, Gribaldo 2011) were
217 updated according to the same procedure. These 73 protein families in combination to the 200
218 newly identified in this study represented 273 phylogenetic markers useful to investigate the
219 evolutionary history of Archaea.

220 We made the choice to not include a bacterial outgroup in our analyses in order to limit long
221 branch attraction (LBA) artifacts which could occur between the long branch leading to the
222 outgroup and the fast evolving lineages within the ingroup, a situation which is frequently
223 encountered when studying ancient evolution (Philippe, Laurent 1998; Gribaldo, Philippe 2002).
224 Moreover, the inclusion of a bacterial outgroup would have severely limited the number of
225 proteins that could be used to study the phylogeny of Archaea.

226 *Supermatrix construction.*

227 The trimmed alignments of the 273 proteins were combined to build five supermatrices.
228 Supermatrices L2, L3 and L4 corresponded to the 200 new proteins (48,904 amino acid
229 positions), the 57 ribosomal proteins (6,228 amino acid positions) and the 16 protein families
230 corresponding to 14 subunits of the RNA polymerase and transcription factors (2,970 amino acid
231 positions), respectively. Supermatrix L1 (35,589 amino acid positions) was built by
232 concatenating the protein families present in more than 119 species (i.e. 106 newly identified
233 proteins, 57 and 16 ribosomal and transcription proteins), which represented less than 10% of
234 missing data per family, whereas the supermatrix XL1 gathered all the 273 proteins (58,102
235 amino acid positions).

236 *Phylogenetic analysis of the supermatrices.*

237 ML trees of the L1, XL1, L2, L3 and L4 amino acid supermatrices were inferred using PhyML
238 version 3.1 (NNI+SPR and tlr options) with the LG+ Γ 8 evolutionary model. The robustness of the
239 resulting trees was assessed using either the non-parametric bootstrap procedure implemented
240 in PhyML (100 replicates of the original alignment) or the SH-like support. Bayesian inferences
241 (BI) were performed using PhyloBayes 3.3b with the CAT+ Γ 8 model. In addition to the amino
242 acid sequence supermatrices, two recoded versions of the L1 supermatrices were analyzed with
243 phylobayes using dayhoff4 and dayhoff6 options (L1-REC4 and L1-REC6, respectively). The
244 four and six Dayhoff's amino acid families corresponded to [(A,G,P,S,T) (D,E,N,Q) (H,K,R)
245 (F,Y,W,I,L,M,V)] plus cysteins treated as missing data (C= ?) and to [(A,G,P,S,T) (D,E,N,Q)
246 (H,K,R) (F,Y,W) (I,L,M,V) (C)]. Two chains were run for at least 10000 cycles, saving one tree in
247 ten. The first 500 trees were discarded as "burn-in" and one on two of the remaining trees from
248 each chain were sampled to test for convergence and to compute the 50% majority rule
249 consensus tree.

250 *Desaturation procedure.*

251 We used two complementary approaches to reduce the mutational saturation level contained in
252 our data.

253 First, a site-by-site desaturation of the L1 supermatrix was carried out with the Slow-Fast method
254 (Brinkmann, Philippe 1999). This site-by-site desaturation strategy was applied only to the L1

255 supermatrix because the larger amount of missing data contained in the XL1 supermatrix could
256 bias the estimation of the evolutionary rates. To do so, we subdivided the sequences of the L1
257 supermatrix into taxonomically balanced groups shown to be monophyletic in previous studies
258 and in this analysis by selecting four to seven sequences per group. The considered groups
259 were Thaumarchaeota + 'Aigarchaeota', Thermoproteales, Sulfolobales, Desulfurococcales,
260 Thermococcales, Methanococcales, Methanobacteriales, Thermoplasmatales + uncultured
261 marine group II + DHVE2, Archaeoglobales, Methanocellales, Methanosarcinales,
262 Methanomicrobiales and Halobacteriales (Supplementary Table S4). Because they were crucial
263 taxa, Korarchaeota and Methanopyrales were also included in the analysis, despite they were
264 each represented by a single species. The evolutionary rate of each amino acid site of the L1
265 supermatrix was estimated with the program SlowFaster (Kostka et al. 2008). 35 alignments (S_0
266 to S_{34}) of increasing size were built by incorporating more and more rapidly-evolving sites.

267 Second, we applied a protein-by-protein desaturation strategy aiming at identifying and focusing
268 the phylogenetic analyses on the slowly evolving proteins. To do so, we estimated the
269 evolutionary rate of each protein by measuring the average evolutionary distances among
270 sequences. Because the taxonomic sampling varied for the different proteins, those average
271 distances among all pairs of sequences were not directly comparable. We thus calculated the
272 average distances among five taxonomic groups widely represented among the 273 protein
273 families: Methanobacteriales, Thermococcales, Sulfolobales, Thermoproteales and
274 Thaumarchaeota + 'Aigarchaeota' to approximate the average evolutionary rate of each marker.
275 We then constructed 28 supermatrices (P_0 to P_{27}) of increasing size by concatenating the protein
276 datasets according to their average evolutionary rates, from the slowest- to the fastest-evolving
277 ones.

278 **Results.**

279 *Expanding the conserved phylogenetic core of Archaea.*

280 Most of the previous works on the phylogeny of Archaea have relied on the analysis of proteins
281 involved in informational systems (in particular ribosomal proteins, subunits of the RNA
282 polymerase, elongation factors, etc.). The in-depth survey of the proteomes of two
283 thaumarchaeotal species *Cenarchaeum symbiosum* and *Nitrosopumilus maritimus* and the
284 aigarchaeon 'Ca. Caldiarchaeum subterraneum' led us to the identification of 200 new
285 conserved proteins widely distributed among Archaea (Figure 1A and Supplementary Table S1).
286 These were either specific to Archaea or strongly supported the monophyly of this domain and
287 its orders, suggesting that they could have been present in the last common ancestor of archaea
288 (LACA) and vertically inherited all along the diversification of this domain.

289 Together with the 57 ribosomal proteins and the 14 RNA polymerase subunits and transcription
290 factors (encoded by 16 genes, because RpoA and RpoB are split in some species) retrieved
291 from previous works (Matte-Tailliez et al. 2002; Brochier, Forterre, Gribaldo 2004), these
292 represented 273 phylogenetic markers of potential interest to investigate the phylogeny of
293 Archaea, most of them being totally new and never used before. Interestingly, the majority of the
294 200 proteins identified by analyzing the proteomes of the two thaumarchaeotes and of 'Ca.
295 Caldiarchaeum subterraneum' did not correspond to genes involved in informational processes

296 (Figure 1B). More precisely, 51 were involved in metabolism, 23 in cellular processes and
297 signaling, 2 in other functions, whereas 47 correspond to poorly characterized protein families.
298 Importantly, to a few exceptions, these protein families were widely distributed in Archaea,
299 indicating that genes involved in non-informational processes can be strongly conserved over
300 large evolutionary times. Beside phylogenetic considerations, the likely presence of these genes
301 in the last common ancestor of Archaea provided interesting biological information on this
302 ancestor. Among the most interesting examples, we detected FtsZ, which, despite being absent
303 in many crenarchaeotes (Bernander 2000), was probably involved in cell division in the archeal
304 ancestor. We also detected enzymes involved in homoserine, thiamine, aspartate, and
305 pseudouridine metabolism, as well as subunits of the exosome and proteasome complexes. The
306 list also includes a large number of ATPase subunits, as well as enzymes involved in
307 hydrogenase expression and maturation, all of them probably related to energy conversion.
308 Given the wide distribution and phyletic patterns of these proteins in archaea, they were most
309 likely present in the last archaeal common ancestor.

310 Archaeal homologues of the whole set of 273 protein families were retrieved from complete
311 genome sequences of 129 species available at the NCBI (Figure 1). These included 34 new
312 species compared to recent studies (Brochier-Armanet, Forterre, Gribaldo 2011), increasing
313 significantly the taxonomic diversity of Archaea that can be considered in massive phylogenomic
314 analyses. The new genomes included the three Nanohaloarchaeota, three additional
315 thaumarchaeotes, six halobacteriales (including the new genera *Natrinema*, *Halopiger*, and
316 *Natronobacterium*), six methanogen Class II and six methanogen Class I (including the new
317 genera *Methanosalsum* and *Methanolinea*), one archaeoglobales, three thermococcales, two
318 desulfurococcales (including the new genus *Pyrolobus*), three sulfobacterales and four
319 thermoproteales (including the new genus *Thermoproteus*).

320 The 273 proteins were checked to identify orthologues and remove paralogous and xenologous
321 sequences. After this trimming step, the alignments were concatenated to generate three
322 supermatrices: L2, L3 and L4 by gathering the 200 newly identified protein families (48,904
323 amino acid positions), the 57 ribosomal proteins (6,228 amino acid positions) and the 16
324 proteins involved in transcription (2,970 amino acid positions), respectively. The size of the
325 ribosomal protein supermatrix L3 was in agreement with previous analyses (Matte-Tailliez et al.
326 2002; Gribaldo, Brochier-Armanet 2006) but smaller than supermatrices from recent reports
327 (Yutin et al. 2012; Podar et al. 2013) despite the use of a similar set of proteins, most likely
328 because of less stringent alignment-trimming criteria in these studies. At this step we removed
329 the seven nanosized archaeal lineages (*Nanoarchaeum equitans*, three Nanohaloarchaeota and
330 three ARMAN species) because of their fast evolutionary rates (Brochier-Armanet, Forterre,
331 Gribaldo 2011; Narasingarao et al. 2012) and their lesser representation in the 273 protein
332 families compared to other species (Figure 1D) in agreement with the relative small size of their
333 genomes, which made them very prone to potential tree reconstruction artifacts (Roure, Baurain,
334 Philippe 2013), meaning that only 122 species were represented in each supermatrix. The ML
335 trees resulting from the three supermatrices showed similar topologies indicating that these
336 supermatrices carried a consistent phylogenetic signal (Supplementary Figures S1-S3). More
337 precisely, the monophyly of the archaeal phyla and orders represented by more than one
338 sequence was recovered and well supported (all SH supports >0.95, except for

339 Desulfurococcales). Moreover, the relationships among the main archaeal orders within each
340 phylum were consistent and in agreement with previous studies (Brochier-Armanet, Forterre,
341 Gribaldo 2011). The strongest discrepancy concerned the position of *Methanopyrus kandleri*,
342 which branched at the base of the Euryarchaeota in the L4 tree, whereas it emerged after the
343 divergence of Thermococcales in the L2 and L3 trees. A similar branching pattern observed in
344 previous analyses of the RNA polymerase and transcription factors was interpreted as the result
345 of a LBA due to the fast evolutionary rate of some components of the *M. kandleri* transcription
346 apparatus (Brochier, Forterre, Gribaldo 2004). The hypothesis of a LBA is strengthened by our
347 new analyses because the basal branching of *M. kandleri* was not observed when we applied
348 Bayesian inference (with the CAT+Γ8 model), less sensible to LBA, to analyze of the same RNA
349 polymerase and transcription factors supermatrix (Supplementary Figures S4).

350 Because they carried a consistent global phylogenetic signal, the 273 protein families were then
351 combined to build two very large supermatrices: L1 (35,589 amino acid positions), obtained by
352 concatenating the 179 proteins present in more than 119 species, and XL1 (58,102 amino acid
353 positions), which resulted from the combination of all the 273 proteins. Without surprise, the ML
354 trees inferred with these two supermatrices (Figure 2) were consistent with the L2, L3 and L4
355 supermatrices (Supplementary Figures S1-S3) and with previous studies (Brochier-Armanet,
356 Forterre, Gribaldo 2011).

357 *Towards a well-resolved global phylogeny of Archaea.*

358 The L1 and the XL1 supermatrices were the largest reported so far to study the phylogeny of
359 Archaea and represented a unique opportunity to decipher in detail the evolutionary history of
360 this domain. We were particularly interested in investigating the quantity and the quality of the
361 phylogenetic signal contained in these protein families and the possible biases which could
362 affect the phylogeny of Archaea. To do so, we applied a two-steps strategy. On the one hand,
363 we applied a site-by-site desaturation strategy, the SlowFast method, to the L1 supermatrix to
364 investigate the relationships among the archaeal phyla and the relative branching pattern of the
365 orders within each phylum. This strategy aimed at reducing the mutational saturation of the data
366 and thus limiting the risk of tree reconstruction artifacts. On the other hand, we used the recoded
367 L1 supermatrices (L1-REC4 and L1-REC6) to reconstruct the BI phylogenies of Euryarchaeota
368 and Crenarchaeota separately to study in detail the evolutionary history of these phyla. This
369 allowed both reducing the mutational saturation and amino acid compositional biases.

370 For the site-by-site desaturation strategy we used the L1 supermatrix to generate 35 SF-
371 supermatrices (S_0 to S_{34}) by including sites with increasing evolutionary rates (Figure 3A) and we
372 inferred ML trees with them Supplementary Figures S5). First of all and without surprise, all the
373 reconstructed trees strongly recovered the monophyly of Crenarchaeota and of Euryarchaeota
374 (Figure 3B). We observed an increasing phylogenetic signal supporting the unrooted
375 [(Thaumarchaeota+'Aigarchaeota',Korarchaeota),(Crenarchaeota,Euryarchaeota)] topology, with
376 the maximal support being associated to the SF-ML-trees S_8 to S_{12} (Figure 3B), but the support
377 in favor of this topology decreased progressively when faster evolving sites were considered.
378 Because we did not include any outgroup sequence to avoid possible LBA artifacts and to
379 maximize the number of conserved alignment positions (see methods), this topology could not
380 be interpreted in terms of evolutionary relationships among the four main archaeal lineages.

381 However, this topology was in agreement with some studies (Elkins et al. 2008), but in
382 contradiction with others (Elkins et al. 2008; Kelly, Wickstead, Gull 2010; Brochier-Armanet,
383 Forterre, Gribaldo 2011; Guy, Ettema 2011; Nunoura et al. 2011; Eme et al. 2013). However,
384 because it was supported by the slowly evolving positions in the site-by-site desaturation
385 approach, this topology could reflect the true structure of the archaeal domain.

386 Furthermore, all the SF-ML trees supported the Thermococcales and Thermoproteales as the
387 first diverging lineages within the Euryarchaeota and Crenarchaeota, respectively (Figure 3B), in
388 agreement with earlier studies (Matte-Tailliez et al. 2002; Gribaldo, Brochier-Armanet 2006;
389 Brochier-Armanet, Forterre, Gribaldo 2011). This provided additional support in favor of the
390 hypothesis that methanogenesis was not the most ancient metabolism, neither in Archaea nor in
391 Euryarchaeota and, therefore, not in Life history (Gribaldo, Brochier-Armanet 2006; Brochier-
392 Armanet, Forterre, Gribaldo 2011), contrarily to what is sometime assumed (Xue et al. 2005;
393 Wong et al. 2007; Yu et al. 2009).

394 The L1-REC6 BI tree of the Crenarchaeota and Euryarchaeota phyla were consistent (Figure 4)
395 and in agreement with previous works (Brochier-Armanet, Forterre, Gribaldo 2011). The
396 recoding of the L1 supermatrix with the Dayhoff4 scheme provided very similar trees (not
397 shown). Regarding Crenarchaeota (Figure 4A), the monophyly of the order Thermoproteales
398 and its genera was clearly supported (all Posterior Probabilities (PP) = 1.0). Similarly, the
399 grouping of the genera *Pyrobaculum* and *Thermoproteus* on the one hand and of *Caldivirga* and
400 *Vulcanisaeta* on the other hand was recovered (PP = 1.0 and PP = 0.82, respectively), whereas
401 *Thermofilum* represented the first lineage diverging within this order (all PP = 1.0). Within
402 Sulfolobales which appeared also as monophyletic (PP = 1.0), the monophyly of *Metallosphaera*
403 was strongly supported as well as its grouping with *Acidianus* (PP = 1.0). In contrast, the
404 monophyly of *Sulfolobus* was not recovered because *Sulfolobus islandicus* and *Sulfolobus*
405 *solfataricus* branched together at the base of the *Acidianus* and *Metallosphaera* group, whereas
406 *Sulfolobus tokodaii* and *Sulfolobus acidocaldarius* represented the first lineage diverging within
407 Sulfolobales (PP = .95). This suggested that the genus *Sulfolobus* may be paraphyletic and
408 encompasses in fact distinct taxonomic lineages.

409 While Thermoproteales and Sulfolobales represented monophyletic orders, the situation was
410 less clear in the case of Desulfurococcales, which were monophyletic in the ML tree of the L1
411 supermatrix (Figure 2) but paraphyletic in the L1-REC6 BI tree, due to the grouping of
412 Sulfolobales with either *Ignicoccus hospitalis* (the host of the nanoarchaeon *Nanoarchaeum*
413 *equitans*) albeit with a weak support (PP = 0.8, Figure 4). The non monophyly of
414 Desulfurococcales was reported and discussed earlier (Brochier-Armanet, Forterre, Gribaldo
415 2011). As for *Sulfolobus*, this could indicate that the order Desulfurococcales is actually non-
416 monophyletic or, conversely, that the relationship between *I. hospitalis* and the Sulfolobales is
417 artefactual. In that case, it could be the result of a few HGT events that were not detected by our
418 trimming procedure or of a LBA artifact between the long stems leading to these two lineages.

419 Regarding the other Desulfurococcales lineages, the monophyly of *Desulfurococcus* and of
420 *Staphylothermus* was strongly supported (PP = 1.0). *Desulfurococcus* grouped robustly with
421 *Thermosphaera aggregans* (PP = 1.0), and represented the sister-lineage of *Staphylothermus*,
422 whereas *Ignisphaera aggregans* emerged more deeply. *Aeropyrum pernix* and *Acidilobus*

423 *saccharovorans* (PP = 1.0) on the one hand and *Hyperthermus butylicus* and *Pyrolobus fumarii*
424 (PP = 1.0) on the other hand, grouped together (PP = 0.93). This confirmed that *Acidilobus* is a
425 true member of Desulfurococcales (Brochier-Armanet, Forterre, Gribaldo 2011) and does not
426 represent a separate order of Crenarchaeota as recently proposed (Prokofeva et al. 2009;
427 Mardanov et al. 2010).

428 As for the Crenarchaeota, the L1 REC6 and L1-REC4 BI euryarchaeal trees were consistent
429 (Figure 4B and not shown) and mainly in agreement with previous work (Brochier-Armanet,
430 Forterre, Gribaldo 2011). The monophyly of each order was recovered and well supported:
431 Thermococcales, Methanobacteriales, Methanococcales, Thermoplasmatales, Halobacteriales,
432 Methanosarcinales, Methanomicrobiales, Methanocellales and Archaeoglobales (all PP = 1.0,
433 Figure 4B). Within Thermococcales, the monophyly of *Pyrococcus* was confirmed (PP = 1.0) but
434 not of *Thermococcus*, due to the early divergence of *Thermococcus barophilus*, *Thermococcus*
435 *sibiricus* and *Thermococcus litoralis* (PP = 1.0), whereas the four remaining *Thermococcus*
436 species formed the sister-lineage of *Pyrococcus* (PP = 1.0). This suggested that the genus
437 *Thermococcus* encompasses distinct lineages. The next diverging lineage gathered
438 methanogens Class I: the Methanobacteriales, Methanopyrales and Methanococcales (PP =
439 1.0), the first two forming a monophyletic group (PP = 0.99), meaning that the monophyly of
440 methanogens Class I was strongly recovered (Figure 4B). These groupings were strengthened
441 by the site-by-site desaturation approach when the slowest evolving sites contained in the SF-
442 supermatrix were considered (Figure 3C, red and blue lines). When faster and faster evolving
443 sites were included (SF-supermatrices S₆ to S₃₄), the support for the monophyly of methanogens
444 Class I decreased progressively and a grouping between Methanococcales and
445 Methanobacteriales appeared. This indicated that the grouping of these two orders and the non
446 monophyly of methanogens Class I observed in the ML L1 and XL1 trees (Figure 2) and in other
447 studies (Brochier-Armanet et al. 2008; Wolf et al. 2012; Yutin et al. 2012; Podar et al. 2013)
448 were likely artefactual. Worth noticing, the grouping of Methanopyrales and Methanobacteriales,
449 which was supported by the slowly evolving sites, was also supported by phylogenies based on
450 genome gene content (Makarova et al. 2007) and by the presence of pseudomurein in their cell
451 wall, a unique biological feature shared by these two archaeal orders (see (Albers, Meyer 2011)
452 and references therein). Pseudomurein could therefore represent a synapomorphy of this
453 lineage.

454 Within Methanobacteriales (Figure 4B), the monophyly of the genera *Methanothermobacter* and
455 *Methanobrevibacter* was well recovered (PP = 1.0), as well as the grouping of *Methanosphaera*
456 and *Methanobacterium* sp. AL-21 (PP = 1.0) together with the genus *Methanobrevibacter* (PP =
457 1.0), whereas *Methanothermobacter* represented the first lineage branching within this order (PP =
458 1.0). In the case of the Methanococcales, the genera *Methanocaldococcus* and *Methanotorris*
459 were monophyletic (PP = 1.0). In contrast, *Methanococcus* appeared paraphyletic due to the
460 clustering of *Methanococcus aeolicus* with *Methanothermococcus okinawensis* (PP = 1.0). This
461 was in agreement with a large scale phylogenomic analysis of Methanococcales based on the
462 universally conserved proteins within this order (Céline Brochier-Armanet and Michel Lecocq,
463 unpublished data) and suggested that *Methanococcus* does not form a monophyletic genus or,
464 conversely, that *Methanothermococcus* should be reclassified within the genus *Methanococcus*.

465 Regarding the intermediate part of the euryarchaeotal tree, the L1-REC6 BI tree and the site-by-
466 site desaturation analysis strongly supported the clustering of Thermoplasmatales, DHEV2
467 (represented by *Aciduliprofundum boonei*) and the uncultured marine group II lineage (PP =
468 1.0), as well as the divergence of Thermoplasmatales and relatives prior the divergence of
469 Archaeoglobales (Figure 2E, and PP = 1.0, Figure 4B). Among these, Thermoplasmatales
470 appeared more closely related to DHEV2 (PP = 1.0) than to the marine group II. However, the
471 very long branch of the latter and its close relationship with Thermoplasmatales in L1-S₁ ML
472 trees inferred with the slowest evolving sites (Figure 3E) deserved consideration: did it result of
473 stochastic effects or could it indicate that basal branching of the marine group II resulted from a
474 LBA? Due to the very restricted taxonomic sampling available for this clade (DHEV2 and marine
475 group II were represented each by only one genome), a definitive answer regarding the
476 phylogenetic position of marine group II would require additional data and analyses. Within
477 Thermoplasmatales, the monophyly of *Thermoplasma* and the sister-relationship between
478 *Ferroplasma* and *Picrophilus* was confirmed (PP = 1.0) (Figure 4B). In contrast, the monophyly of
479 the genus *Archaeoglobus* within the Archaeoglobales was not recovered, due to the basal
480 branching of *Archaeoglobus profundus*, albeit with a weak support (PP = 0.83). The apical part
481 of the euryarchaeotal tree gathered Halobacteriales and the Methanomicrobia, represented by
482 three orders of methanogens: Methanomicrobiales, Methanocellales and Methanosarcinales.
483 These orders were proposed to form a second class of methanogens, tentatively called
484 methanogens Class II, representing the sister-lineage of the Halobacteriales (Baptiste, Brochier,
485 Boucher 2005). While the monophyly of methanogens Class II was recovered in some
486 phylogenetic analyses (Wolf et al. 2012; Yutin et al. 2012), it was not in others (Brochier-
487 Armanet et al. 2008; Csuros, Miklos 2009; Kelly, Wickstead, Gull 2010; Groussin, Gouy 2011;
488 Podar et al. 2013), and thus remained debated (Brochier-Armanet, Forterre, Gribaldo 2011). In
489 our study, the monophyly of methanogens Class II was recovered neither in the L1-REC6 BI tree
490 nor by the site-by-site desaturation strategy (Figure 4F and Figure 3B). The L1-REC6 BI tree
491 strongly supported the emergence of Halobacteriales within Methanomicrobia as the sister-
492 lineage of Methanomicrobiales (PP = 1.0), and the early divergence of Methanosarcinales (PP =
493 1.0). However, the picture provided by the analysis of the slowest evolving positions was slightly
494 more complex and highlighted potential artifacts. First, the SF-ML trees based on the slowest
495 evolving positions (S₁ to S₆, Figure 3F) favored a sister-grouping between Methanocellales and
496 Halobacteriales, with Methanosarcinales as the first diverging lineage within methanogens Class
497 II (S₁ to S₄, Figure 3D), suggesting that the branching pattern observed in ML and BI trees based
498 on the complete L1 supermatrix could be artefactual (Figure 2B and Figure 4B). The
499 relationships among methanogens Class II and Halobacteriales shifted dramatically following the
500 addition of 1376 positions (transition from S₆ to S₇): Methanocellales moved to the base of
501 methanogens Class II, whereas Halobacteriales grouped with Methanomicrobiales (Figures 2D
502 and 2F). Interestingly, while the association between Halobacteriales and Methanomicrobiales
503 remained stable in the ML trees built with SF-supermatrices with an index greater than 20
504 (Figure 2F), an association between Methanosarcinales and Methanocellales appeared (Figure
505 2D). However, this relation was highly suspicious because it occurred only when the fastest
506 evolving positions of the L1 supermatrix were taken into account.

507

508 *Origin of the nanosized archaeal lineages.*

509 One of the most intriguing questions regarding the phylogeny of Archaea regarded the
510 phylogenetic position of the newly discovered nanosized lineages, namely the two uncultured
511 ARMAN lineages (Baker et al. 2006), the Nanoarchaeota (Huber et al. 2002) and the uncultured
512 Nanohaloarchaea (Narasingarao et al. 2012). Due to fast evolutionary rates, their phylogenetic
513 position has been particularly difficult to assess and is debated (Brochier-Armanet, Forterre,
514 Gribaldo 2011; Podar et al. 2013). To address this question we applied a protein-by-protein
515 desaturation strategy aiming at identifying the slowly evolving proteins to build 28 supermatrices
516 (P_0 to P_{27}) by including faster and faster evolving proteins. The supermatrices gathering the
517 slowest evolving proteins were expected to contain the most reliable phylogenetic signal. The
518 ML trees based on the 10 slowly evolving genes (P_0 , not shown) confirmed the sister-
519 relationship between Nanohaloarchaea and Halobacteria (SH = 0.998) (Narasingarao et al.
520 2012). It also strongly supported the late divergence of '*Ca. Micrarchaeum acidiphilum*' within
521 Euryarchaeota in agreement with SSU rRNA analyses (Baker et al. 2006; Baker et al. 2010) and
522 a recent report (Brochier-Armanet, Forterre, Gribaldo 2011), but in contradiction analyses based
523 on single protein (RadA, EF-2, EF-1A and SecY) (Baker et al. 2010). More precisely '*Ca.*
524 '*Micrarchaeum acidiphilum*' branched robustly at the base of Thermoplasmatales, DHVE2 and
525 group II (SH = 0.977). Finally, '*Ca. Parvarchaeum acidophilus*' and *Nanoarchaeum equitans*
526 grouped together (SH = 0.925) at the base of the Thermococcales (SH = 0.493) consistently with
527 a recent report (Brochier-Armanet, Forterre, Gribaldo 2011). However, this relationship should
528 be considered with caution because it could result from a LBA given the very long branches of
529 these two nanosized lineages. Strengthening the hypothesis of a LBA, the two '*Ca.*
530 '*Parvarchaeum*' branched robustly at the base of the Thermoplasmatales, DHVE2 and group II
531 (SH = 0.997) in the P2 ML tree, when the other fast evolving lineages were removed from the
532 analysis (not shown).

533 Finally, we have showed that the L1 supermatrix contained a reliable phylogenetic signal,
534 especially when the slowly evolving positions were considered. From our previous analysis, we
535 established that the best ratio between phylogenetic signal and noise was found in the L1-S₂ to
536 L1-S₅ supermatrices, and that a shift occurred between the L1-S₆ and L1-S₇ supermatrices.
537 Accordingly, we used these supermatrices to test the phylogenetic position of *Nanoarchaeum*
538 *equitans*. To do so, added the sequences of this tiny archaea in the S₁ to S₆ supermatrices.
539 Their phylogenetic analysis provided trees supporting the grouping of Nanoarchaeota within
540 Euryarchaeota and more precisely as the sister-lineage of Thermococcales, albeit with moderate
541 supports, whereas a branching at the base of Euryarchaeota as observed for the S₇ and S₈
542 matrices (not shown).

543 **Discussion**

544 Taking the opportunity of the recent burst of complete genome sequences covering an ever-
545 growing part of the archaeal diversity we have identified 200 proteins of phylogenetic interest
546 which can be used in combination with the 73 proteins of the ribosome (Matte-Tailliez et al.
547 2002) and transcription apparatus (Brochier, Forterre, Gribaldo 2004) that have been employed
548 in recent years to investigate the phylogeny of the third domain of Life. The discovery that the
549 transcription and the translation apparatus carry a consistent phylogenetic signal was an

550 important step because it highlighted the existence of a conserved phylogenetic core of genes
551 useful to study the ancient evolution of Archaea (Brochier, Forterre, Gribaldo 2005). However,
552 one could argue that transcription and translation do not represent independent systems
553 because they are coupled in archaea (French et al. 2007), as it is also the case in bacteria.
554 Therefore, it can be speculated that they have coevolved (including gene losses and horizontal
555 gene transfers), explaining why they carry a similar phylogenetic signal. Most of the 200 markers
556 identified in this study are not functionally linked to translation or transcription. More important,
557 the majority are not involved in informational processes. Therefore, the consistent phylogenetic
558 signal carried by these markers and the transcription and translation apparatuses cannot be
559 explained by functional links. This is a strong argument for the existence of a conserved core of
560 genes carrying a phylogenetic signal that reflects the evolutionary relationships of organisms.
561 This contrasts with recent studies claiming that horizontal gene transfers have overwhelmed the
562 evolutionary history of prokaryotes, which cannot be represented by a tree and should be
563 replaced by a network (or a rhizome, etc.) (Dagan, Martin 2006). However, networks of genes
564 can neither capture the complexity of the evolutionary history of living beings (Gribaldo, Brochier
565 2009). In particular, they are not suitable to represent speciation events or gene duplications and
566 losses that are important processes of biological evolution. On the contrary, deciphering the
567 vertical relationships among organisms provides the indispensable framework to analyze the
568 evolutionary history of the genomes (including gene transfers, duplications and losses).
569 Altogether, the 273 genes composing the conserved phylogenetic core of Archaea represent 10-
570 20% of the gene content in typical archaeal genomes. This is one order of magnitude more than
571 the provocative picture of the tree of 1% (Dagan, Martin 2006). This conserved phylogenetic
572 core of genes, which were inherited from the last common ancestor of Archaea, is not strictly
573 conserved in all lineages, meaning that differential gene losses or punctual horizontal gene
574 transfers could have occurred in one or the other lineage during the diversification of this
575 domain. It corresponds to the “soft” conserved phylogenetic core of genes we defined in 2006
576 (Gribaldo, Brochier-Armanet 2006), which is much larger than the strictly conserved
577 phylogenetic core of genes that is the set of universal archaeal genes. While the latter is fated
578 to disappear when a larger and larger part of the diversity will be taken into account in the
579 analyses, the former is not. On the contrary, taking into consideration more genomes and more
580 lineages allows distinguishing accessory genes composing the variable shell from those which
581 were ancestrally present and punctually lost in each lineage, thus increasing the size of the soft
582 conserved phylogenetic core of genes. For instance, taking into account ‘*Ca. Caldiarchaeum*
583 *subterraneum*’ (the sole representative of the candidate phylum ‘Aigarchaeota’) allowed
584 identifying 43 additional protein families that were not detected when we used as reference the
585 genomes of the two thaumarchaeota *Cenarchaeum symbiosum* and *Nitrosopumilus maritimus*.

586 The accurate analysis of the 273 markers through supermatrix and desaturation approaches
587 brought a new light on the evolutionary history of Archaea. First, we confirmed a number of
588 nodes that were proposed before (e.g. the early emergence of Thermococcales within
589 Euryarchaeota, the grouping of Thermoplasmatales with DHEV2 and the uncultured marine
590 group II, or the branching of the candidate phylum ‘Aigarchaeota’ at the base of
591 Thaumarchaeota). We also pointed some discrepancies between the phylogeny of Archaea and
592 the current taxonomy of this domain (e.g. the refutation of the order Acidilobales, or the
593 paraphyly of the genera *Methanococcus* and *Sulfolobus*) that deserves further investigation and

594 could lead to revise the current archaeal taxonomy. More importantly, we settled a number of
595 important debated relationships such as the reliability of the methanogens Class I which
596 appeared strongly supported by the slowest evolving sites and by Bayesian inferences.
597 Therefore, it represents a *bona fide* class for which we propose the name Methanohoma (from
598 the Greek ὁμάς meaning group), accordingly the classes Methanococci/Methanotherma,
599 Methanobacteria and Methanopyri should be considered as sub-classes. The apparent non-
600 monophyly of this class observed in some phylogenetic analyses was most likely due to tree
601 reconstruction artifacts and/or biases. We also provided some support for the monophyly of
602 Desulfurococcales. While significant advances emerged from our analyses, a few relationships
603 remained elusive. The most important concerned the relationship among the methanogens
604 Class II and Halobacteriales, which can reflect a lack of phylogenetic signal despite the use of
605 very large supermatrices and accurate desaturation approaches and/or biases linked to HGT
606 that would have escaped our trimming procedure. Even if a close relationship between
607 Methanocellales and Halobacteriales seemed to emerge, additional data and analyses will be
608 required to definitively close the question. Finally our site-by-site and protein-by-protein
609 desaturation analyses clarified the origin of nanosized archaeal lineages. We showed that they
610 are bona fide Euryarchaeota. More precisely, we confirmed the sister-ship between
611 Nanohaloarchaea and Halobacteria. We refined the phylogenetic position of 'Ca. Micrarchaeum
612 acidiphilum' and supported their emergence of at the base of the large group comprising
613 Thermoplasmatales, DHVE2 and group II. In contrast, additional studies are required to
614 definitively conclude about the position of the two 'Ca. Parvarchaeum'. The analysis of the slowly
615 evolving positions and proteins both confirmed the relationship between Thermococcales and
616 Nanoarchaeota proposed nearly ten year ago (Brochier et al. 2005) and was in agreement with a
617 recent analysis (Brochier-Armanet, Forterre, Gribaldo 2011). This suggested that their
618 emergence at the base of Euryarchaeota recently reported (Podar et al. 2013) was the
619 consequence of LBA. LBA could also explain a number of conflicting relationships observed.
620 This was the case for instance of the grouping of Methanococcales and Methanobacteriales or
621 of the Halobacteriales and Methanomicrobiales. The use of desaturation approaches, and
622 accurate reconstruction methods and evolutionary models can overcome these artifacts. The
623 use of rough phylogenetic reconstruction methods (e.g. distance or approximate maximum
624 likelihood ones) or non-realistic evolutionary models should thus be avoided for phylogenomic
625 studies aiming at reconstructing the very ancient relationships.

626

627 **Figure legends.**

628 **Figure 1.** Number of archaeal orthologues **(A)**, functional distribution **(B)** and **(D)** taxonomic
629 distribution of the 200 newly identified protein markers. **(C)** Number of archaeal orthologues of
630 transcription (grey) and ribosomal proteins (black).

631 **Figure 2.** Phylogeny of Archaea inferred with the L1 supermatrix (122 species, 35,589 amino
632 acid positions). This provisional tree was inferred with FastTree (WAG+I8). The scale bars
633 represent the average number of substitutions per site. Numbers at branches represent ML SH-
634 like supports estimated with the L1 and XL1 supermatrices, respectively.

635 **Figure 3. (A)** Number of positions considered in each SF-matrix (S_0 to S_{34}). **(B-F)** Evolution of
636 the SH-supports for major archaeal relationships according to the site-by-site desaturation
637 strategy. X and Y axes correspond to the SF-matrices and to the SH-supports, respectively (see
638 legends on the Figure).

639 **Figure 4.** BI trees of Crenarchaeota (33 species) **(A)** and Euryarchaeota (79 species) **(B)**, based
640 on the L1-REC6 supermatrix (35,589 positions). The tree was inferred with Phylobayes
641 (CAT+ Γ 8+Dayhoff6). The scale bars represent the average number of substitutions per site.
642 Numbers at branches are posterior probabilities.

643

644 **Supplementary material.**

645 **Supplementary Table S1.** Table showing the taxonomic sampling used for the identification of
646 new markers that can be potentially used to study the phylogeny of Archaea.

647 **Supplementary Table S2.** Table showing the 200 new proteins used for the phylogenetic
648 analysis.

649 **Supplementary Table S3.** Table showing the taxonomic distribution of the 129 archaeal
650 genomes used for the assembly of final datasets.

651 **Supplementary Figure S1.** Unrooted maximum likelihood phylogeny of Archaea based on the
652 200 newly identified proteins (L2 supermatrix, 48,904 amino acid positions, 122 species). The
653 tree was inferred with PhyML (LG+ Γ 8). The scale bar represents the average number of
654 substitutions per site. Numbers at branches are SH-like supports

655 **Supplementary Figure S2.** Unrooted maximum likelihood phylogeny of Archaea based on the
656 53 ribosomal proteins (L3 supermatrix, 6,228 amino acid positions, 122 organisms). The tree
657 was inferred with PhyML (LG+ Γ 8). The scale bar represents the average number of substitutions
658 per site. Numbers at branches are SH-like supports.

659 **Supplementary Figure S3.** Unrooted maximum likelihood phylogeny of Archaea based on the
660 15 subunits of the RNA polymerase and transcription factors (L4 supermatrix, 2,970 amino acid
661 positions, 122 species). The tree was inferred with PhyML (LG+ Γ 8). The scale bar represents
662 the average number of substitutions per site. Numbers at branches are SH-like supports.

663 **Supplementary Figure S4.** Unrooted Bayesian tree of Archaea based on the 15 subunits of the
664 RNA polymerase and transcription factors (L4 supermatrix, 2970 amino acid positions, 122
665 species). The tree was inferred with PhyloBayes (CAT+ Γ 8). The scale bar represents the
666 average number of substitutions per site. Numbers at branches are posterior probabilities.

667 **Supplementary Figure S5.** Unrooted maximum likelihood trees resulting from the site-by-site
668 desaturation strategy. The S_0 to S_{34} supermatrices were extracted from the XL1 supermatrix (see
669 material and methods). Trees were inferred with PhyML (CAT+ Γ 8). The scale bar represents the
670 average number of substitutions per site. Numbers at branches are SH supports.

671

672 **Acknowledgments.**

673 This work was supported by the French National Agency for Research (EVOLDEEP project,
674 contract number ANR-08-GENM-024-002) and by the Investissement d'Avenir grant (ANR-10-
675 BINF-01-01). C.B-A is member of the Institut Universitaire de France. C. Petitjean was the
676 recipient of a grant from the ANR EvolDeep.

677

678

679 **Bibliography.**

680

- 681 Albers, SV, BH Meyer. 2011. The archaeal cell envelope. *Nat Rev Microbiol* 9:414-426.
- 682 Altschul, SF, TL Madden, AA Schaffer, J Zhang, Z Zhang, W Miller, DJ Lipman. 1997. Gapped BLAST
683 and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res*
684 25:3389-3402.
- 685 Baker, BJ, LR Comolli, GJ Dick, LJ Hauser, D Hyatt, BD Dill, ML Land, NC Verberkmoes, RL Hettich, JF
686 Banfield. 2010. Enigmatic, ultrasmall, uncultivated Archaea. *Proc Natl Acad Sci U S A* 107:8806-
687 8811.
- 688 Baker, BJ, GW Tyson, RI Webb, J Flanagan, P Hugenholtz, EE Allen, JF Banfield. 2006. Lineages of
689 acidophilic archaea revealed by community genomic analysis. *Science* 314:1933-1935.
- 690 Bapteste, E, C Brochier, Y Boucher. 2005. Higher-level classification of the Archaea: evolution of
691 methanogenesis and methanogens. *Archaea* 1:353-363.
- 692 Bernander, R. 2000. Chromosome replication, nucleoid segregation and cell division in archaea. *Trends*
693 *Microbiol* 8:278-283.
- 694 Bodilis, J, S Nsique-Meilo, L Besaury, L Quillet. 2012. Variable copy number, intra-genomic
695 heterogeneities and lateral transfers of the 16S rRNA gene in *Pseudomonas*. *PLoS One*
696 7:e35647.
- 697 Brinkmann, H, H Philippe. 1999. Archaea sister group of Bacteria? Indications from tree reconstruction
698 artifacts in ancient phylogenies. *Mol Biol Evol* 16:817-825.
- 699 Brochier-Armanet, C, B Boussau, S Gribaldo, P Forterre. 2008. Mesophilic Crenarchaeota: proposal for a
700 third archaeal phylum, the Thaumarchaeota. *Nat Rev Microbiol* 6:245-252.
- 701 Brochier-Armanet, C, P Forterre, S Gribaldo. 2011. Phylogeny and evolution of the Archaea: one hundred
702 genomes later. *Curr Opin Microbiol* 14:274-281.
- 703 Brochier, C, P Forterre, S Gribaldo. 2004. Archaeal phylogeny based on proteins of the transcription and
704 translation machineries: tackling the *Methanopyrus kandleri* paradox. *Genome Biol* 5:R17.
- 705 Brochier, C, P Forterre, S Gribaldo. 2005. An emerging phylogenetic core of Archaea: phylogenies of
706 transcription and translation machineries converge following addition of new genome sequences.
707 *BMC Evol Biol* 5:36.
- 708 Brochier, C, S Gribaldo, Y Zivanovic, F Confalonieri, P Forterre. 2005. Nanoarchaea: representatives of a
709 novel archaeal phylum or a fast-evolving euryarchaeal lineage related to Thermococcales?
710 *Genome Biol* 6:R42.
- 711 Brochier, C, H Philippe. 2002. A non-hyperthermophilic ancestor for bacteria. *Nature* 417:244.
- 712 Cavicchioli, R. 2011. Archaea--timeline of the third domain. *Nat Rev Microbiol* 9:51-61.
- 713 Criscuolo, A, S Gribaldo. 2010. BMGE (Block Mapping and Gathering with Entropy): a new software for
714 selection of phylogenetic informative regions from multiple sequence alignments. *BMC Evol Biol*
715 10:210.
- 716 Csuros, M, I Miklos. 2009. Streamlining and large ancestral genomes in Archaea inferred with a
717 phylogenetic birth-and-death model. *Mol Biol Evol* 26:2087-2095.
- 718 Dagan, T, W Martin. 2006. The tree of one percent. *Genome Biol* 7:118.
- 719 Delsuc, F, H Brinkmann, H Philippe. 2005. Phylogenomics and the reconstruction of the tree of life. *Nat*
720 *Rev Genet* 6:361-375.
- 721 Elkins, JG, M Podar, DE Graham, et al. 2008. A korarchaeal genome reveals insights into the evolution of
722 the Archaea. *Proc Natl Acad Sci U S A* 105:8102-8107.
- 723 Eme, L, LJ Reigstad, A Spang, A Lanzen, T Weinmaier, T Rattei, C Schleper, C Brochier-Armanet. 2013.
724 Metagenomics of Kamchatkan hot spring filaments reveal two new major (hyper)thermophilic
725 lineages related to Thaumarchaeota. *Res Microbiol* 164:425-438.
- 726 Felsenstein, J. 1978. Cases in which parsimony or compatibility methods will be positively misleading.
727 *Syst Zool* 27:401-410.
- 728 French, SL, TJ Santangelo, AL Beyer, JN Reeve. 2007. Transcription and translation are coupled in
729 Archaea. *Mol Biol Evol* 24:893-895.

730 Galtier, N, JR Lobry. 1997. Relationships between genomic G+C content, RNA secondary structures, and
731 optimal growth temperature in prokaryotes. *J Mol Evol* 44:632-636.

732 Gribaldo, S, C Brochier-Armanet. 2006. The origin and evolution of Archaea: a state of the art. *Philos*
733 *Trans R Soc Lond B Biol Sci* 361:1007-1022.

734 Gribaldo, S, C Brochier. 2009. Phylogeny of prokaryotes: does it exist and why should we care? *Res*
735 *Microbiol* 160:513-521.

736 Gribaldo, S, H Philippe. 2002. Ancient phylogenetic relationships. *Theor Popul Biol* 61:391-408.

737 Groussin, M, B Boussau, M Gouy. 2013. A Branch-Heterogeneous Model of Protein Evolution for Efficient
738 Inference of Ancestral Sequences. *Syst Biol*.

739 Groussin, M, M Gouy. 2011. Adaptation to environmental temperature is a major determinant of molecular
740 evolutionary rates in archaea. *Mol Biol Evol* 28:2661-2674.

741 Guindon, S, JF Dufayard, V Lefort, M Anisimova, W Hordijk, O Gascuel. 2010. New algorithms and
742 methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0.
743 *Syst Biol* 59:307-321.

744 Guy, L, TJ Ettema. 2011. The archaeal 'TACK' superphylum and the origin of eukaryotes. *Trends*
745 *Microbiol* 19:580-587.

746 Hallam, SJ, KT Konstantinidis, N Putnam, C Schleper, Y Watanabe, J Sugahara, C Preston, J de la Torre,
747 PM Richardson, EF DeLong. 2006. Genomic analysis of the uncultivated marine crenarchaeote
748 *Cenarchaeum symbiosum*. *Proc Natl Acad Sci U S A* 103:18296-18301.

749 Huber, H, MJ Hohn, R Rachel, T Fuchs, VC Wimmer, KO Stetter. 2002. A new phylum of Archaea
750 represented by a nanosized hyperthermophilic symbiont. *Nature* 417:63-67.

751 Katoh, K, DM Standley. 2013. MAFFT multiple sequence alignment software version 7: improvements in
752 performance and usability. *Mol Biol Evol* 30:772-780.

753 Kelly, S, B Wickstead, K Gull. 2010. Archaeal phylogenomics provides evidence in support of a
754 methanogenic origin of the Archaea and a thaumarchaeal origin for the eukaryotes. *Proc Biol Sci*.

755 Kitahara, K, K Miyazaki. 2013. Revisiting bacterial phylogeny: Natural and experimental evidence for
756 horizontal gene transfer of 16S rRNA. *Mob Genet Elements* 3:e24210.

757 Kostka, M, M Uzlikova, I Cepicka, J Flegr. 2008. SlowFaster, a user-friendly program for slow-fast analysis
758 and its application on phylogeny of *Blastocystis*. *BMC Bioinformatics* 9:341.

759 Lartillot, N, H Brinkmann, H Philippe. 2007. Suppression of long-branch attraction artefacts in the animal
760 phylogeny using a site-heterogeneous model. *BMC Evol Biol* 7 Suppl 1:S4.

761 Lartillot, N, H Philippe. 2004. A Bayesian mixture model for across-site heterogeneities in the amino-acid
762 replacement process. *Mol Biol Evol* 21:1095-1109.

763 Le, SQ, O Gascuel. 2008. An improved general amino acid replacement matrix. *Mol Biol Evol* 25:1307-
764 1320.

765 Le, SQ, N Lartillot, O Gascuel. 2008. Phylogenetic mixture models for proteins. *Philos Trans R Soc Lond*
766 *B Biol Sci* 363:3965-3976.

767 Lopez-Garcia, P, D Moreira. 2008. Tracking microbial biodiversity through molecular and genomic
768 ecology. *Res Microbiol* 159:67-73.

769 Makarova, KS, AV Sorokin, PS Novichkov, YI Wolf, EV Koonin. 2007. Clusters of orthologous genes for
770 41 archaeal genomes and implications for evolutionary genomics of archaea. *Biol Direct* 2:33.

771 Mardanov, AV, VA Svetlitchnyi, AV Beletsky, MI Prokofeva, EA Bonch-Osmolovskaya, NV Ravin, KG
772 Skryabin. 2010. The genome sequence of the crenarchaeon *Acidilobus saccharovorans* supports
773 a new order, Acidilobales, and suggests an important ecological role in terrestrial acidic hot
774 springs. *Appl Environ Microbiol* 76:5652-5657.

775 Matte-Tailliez, O, C Brochier, P Forterre, H Philippe. 2002. Archaeal phylogeny based on ribosomal
776 proteins. *Mol Biol Evol* 19:631-639.

777 Narasingarao, P, S Podell, JA Ugalde, C Brochier-Armanet, JB Emerson, JJ Brocks, KB Heidelberg, JF
778 Banfield, EE Allen. 2012. De novo metagenomic assembly reveals abundant novel major lineage
779 of Archaea in hypersaline microbial communities. *ISME J* 6:81-93.

780 Nunoura, T, Y Takaki, J Kakuta, et al. 2011. Insights into the evolution of Archaea and eukaryotic protein
781 modifier systems revealed by the genome of a novel archaeal group. *Nucleic Acids Res* 39:3204-
782 3223.

783 Philippe, H, A Germot. 2000. Phylogeny of eukaryotes based on ribosomal RNA: long-branch attraction
784 and models of sequence evolution. *Mol Biol Evol* 17:830-834.

785 Philippe, H, J Laurent. 1998. How good are deep phylogenetic trees? *Curr Opin Genet Dev* 8:616-623.

786 Philippe, H, P Lopez, H Brinkmann, K Budin, A Germot, J Laurent, D Moreira, M Muller, H Le Guyader.
787 2000. Early-branching or fast-evolving eukaryotes? An answer based on slowly evolving positions.
788 Proc Biol Sci 267:1213-1221.

789 Podar, M, KS Makarova, DE Graham, YI Wolf, EV Koonin, AL Reysenbach. 2013. Insights into archaeal
790 evolution and symbiosis from the genomes of a nanoarchaeon and its inferred crenarchaeal host
791 from Obsidian Pool, Yellowstone National Park. Biol Direct 8:9.

792 Price, MN, PS Dehal, AP Arkin. 2010. FastTree 2--approximately maximum-likelihood trees for large
793 alignments. PLoS One 5:e9490.

794 Prokofeva, MI, NA Kostrikina, TV Kolganova, TP Tourova, AM Lysenko, AV Lebedinsky, EA Bonch-
795 Osmolovskaya. 2009. Isolation of the anaerobic thermoacidophilic crenarchaeote *Acidilobus*
796 *saccharovorans* sp. nov. and proposal of *Acidilobales* ord. nov., including *Acidilobaceae* fam. nov.
797 and *Caldisphaeraceae* fam. nov. Int J Syst Evol Microbiol 59:3116-3122.

798 Rinke, C, P Schwientek, A Sczyrba, et al. 2013. Insights into the phylogeny and coding potential of
799 microbial dark matter. Nature.

800 Robertson, CE, JK Harris, JR Spear, NR Pace. 2005. Phylogenetic diversity and ecology of environmental
801 Archaea. Curr Opin Microbiol 8:638-642.

802 Rodriguez-Ezpeleta, N, H Brinkmann, G Burger, AJ Roger, MW Gray, H Philippe, BF Lang. 2007. Toward
803 resolving the eukaryotic tree: the phylogenetic positions of jakobids and cercozoans. Curr Biol
804 17:1420-1425.

805 Roger, AJ. 1999. Reconstructing early events in eukaryotic evolution. Am. Nat. 154:S146-S163.

806 Roure, B, D Baurain, H Philippe. 2013. Impact of missing data on phylogenies inferred from empirical
807 phylogenomic data sets. Mol Biol Evol 30:197-214.

808 Stiller, J, B Hall. 1999. Long-branch attraction and the rDNA model of early eukaryotic evolution. Mol Biol
809 Evol 16:1270-1279.

810 Walker, CB, JR de la Torre, MG Klotz, et al. 2010. Nitrosopumilus maritimus genome reveals unique
811 mechanisms for nitrification and autotrophy in globally distributed marine crenarchaea. Proc Natl
812 Acad Sci U S A 107:8818-8823.

813 Waters, E, MJ Hohn, I Ahel, et al. 2003. The genome of *Nanoarchaeum equitans*: insights into early
814 archaeal evolution and derived parasitism. Proc Natl Acad Sci U S A 100:12984-12988.

815 Woese, CR, L Achenbach, P Rouviere, L Mandelco. 1991. Archaeal phylogeny: reexamination of the
816 phylogenetic position of *Archaeoglobus fulgidus* in light of certain composition-induced artifacts.
817 Syst Appl Microbiol 14:364-371.

818 Woese, CR, GE Fox. 1977. Phylogenetic structure of the prokaryotic domain: the primary kingdoms. Proc.
819 Natl. Acad. Sci. USA 74:5088-5090.

820 Wolf, YI, KS Makarova, N Yutin, EV Koonin. 2012. Updated clusters of orthologous genes for Archaea: a
821 complex ancestor of the Archaea and the byways of horizontal gene transfer. Biol Direct 7:46.

822 Wong, JT, J Chen, WK Mat, SK Ng, H Xue. 2007. Polyphasic evidence delineating the root of life and
823 roots of biological domains. Gene 403:39-52.

824 Wu, D, P Hugenholtz, K Mavromatis, et al. 2009. A phylogeny-driven genomic encyclopaedia of Bacteria
825 and Archaea. Nature 462:1056-1060.

826 Xue, H, SK Ng, KL Tong, JT Wong. 2005. Congruence of evidence for a *Methanopyrus*-proximal root of
827 life based on transfer RNA and aminoacyl-tRNA synthetase genes. Gene 360:120-130.

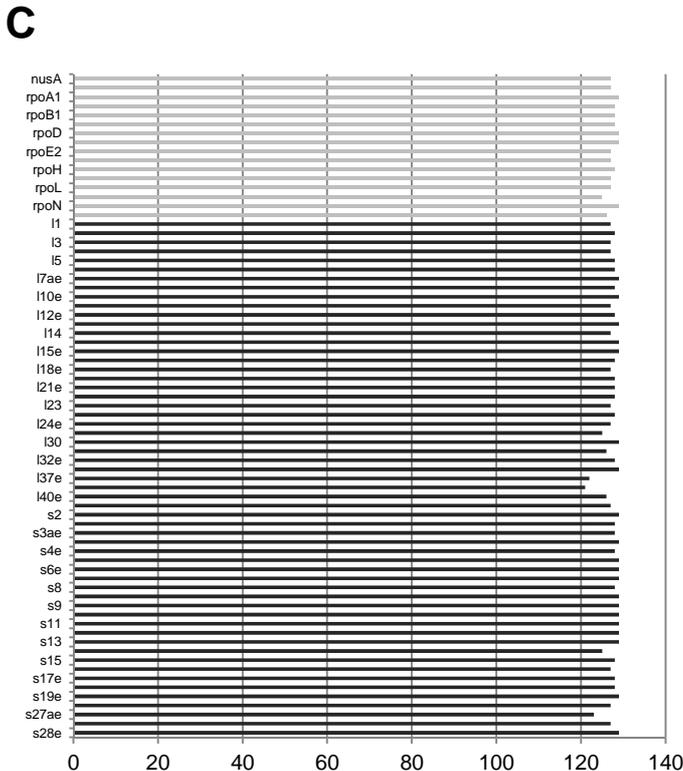
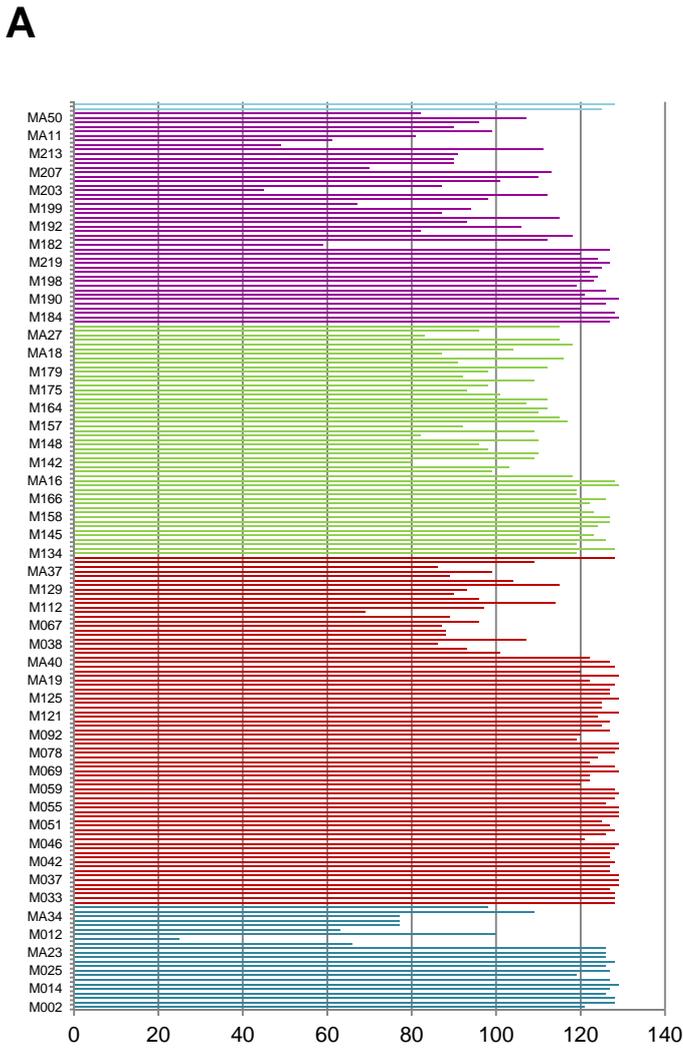
828 Yap, WH, Z Zhang, Y Wang. 1999. Distinct types of rRNA operons exist in the genome of the
829 actinomycete *Thermomonospora chromogena* and evidence for horizontal transfer of an entire
830 rRNA operon. J Bacteriol 181:5201-5209.

831 Yu, Z, K Takai, A Slesarev, H Xue, JT Wong. 2009. Search for primitive *Methanopyrus* based on genetic
832 distance between Val- and Ile-tRNA synthetases. J Mol Evol 69:386-394.

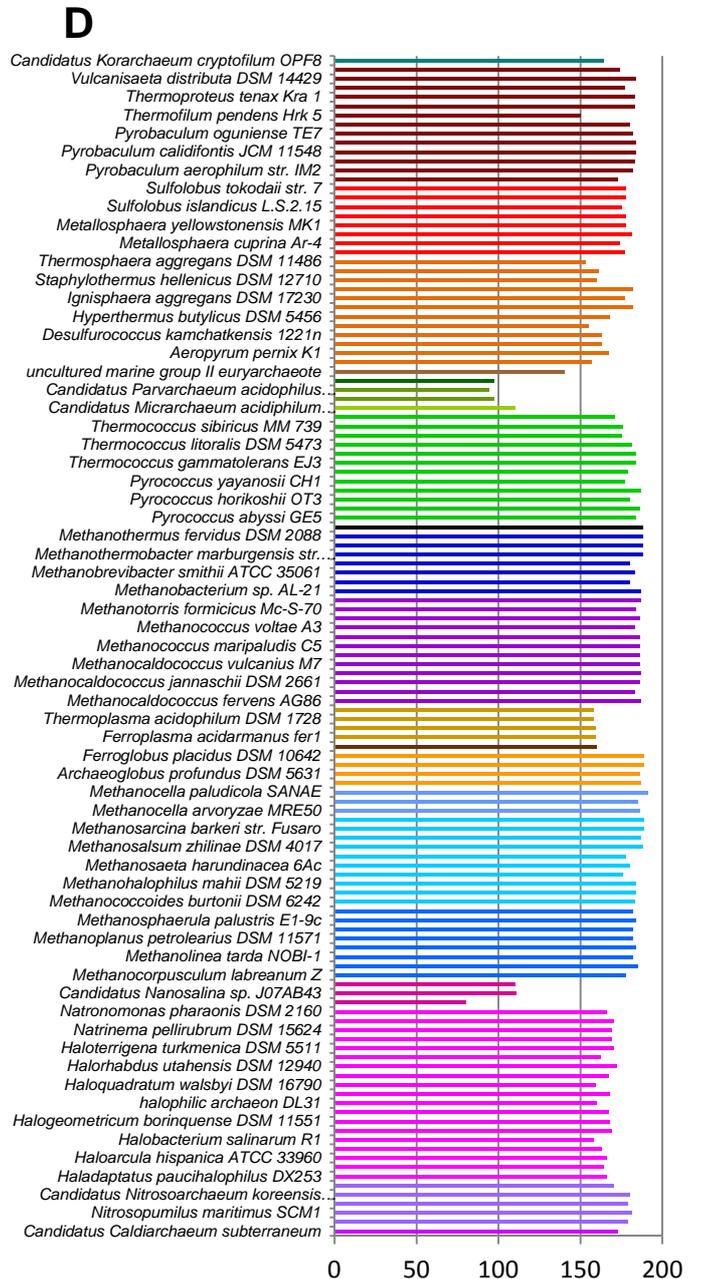
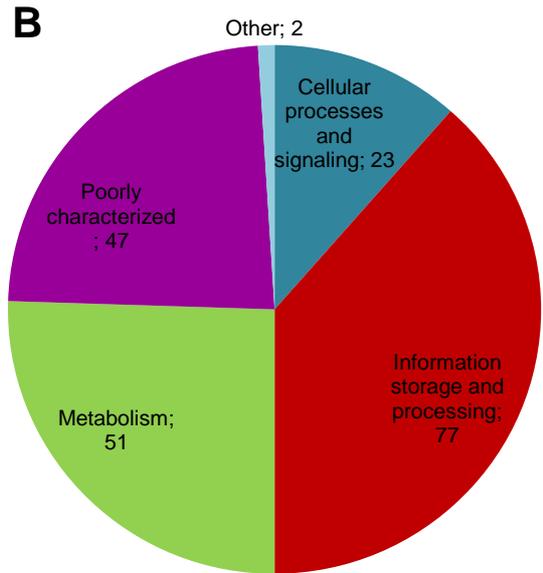
833 Yutin, N, P Puigbò, EV Koonin, YI Wolf. 2012. Phylogenomics of Prokaryotic Ribosomal Proteins. PLoS
834 One 7:e36972.

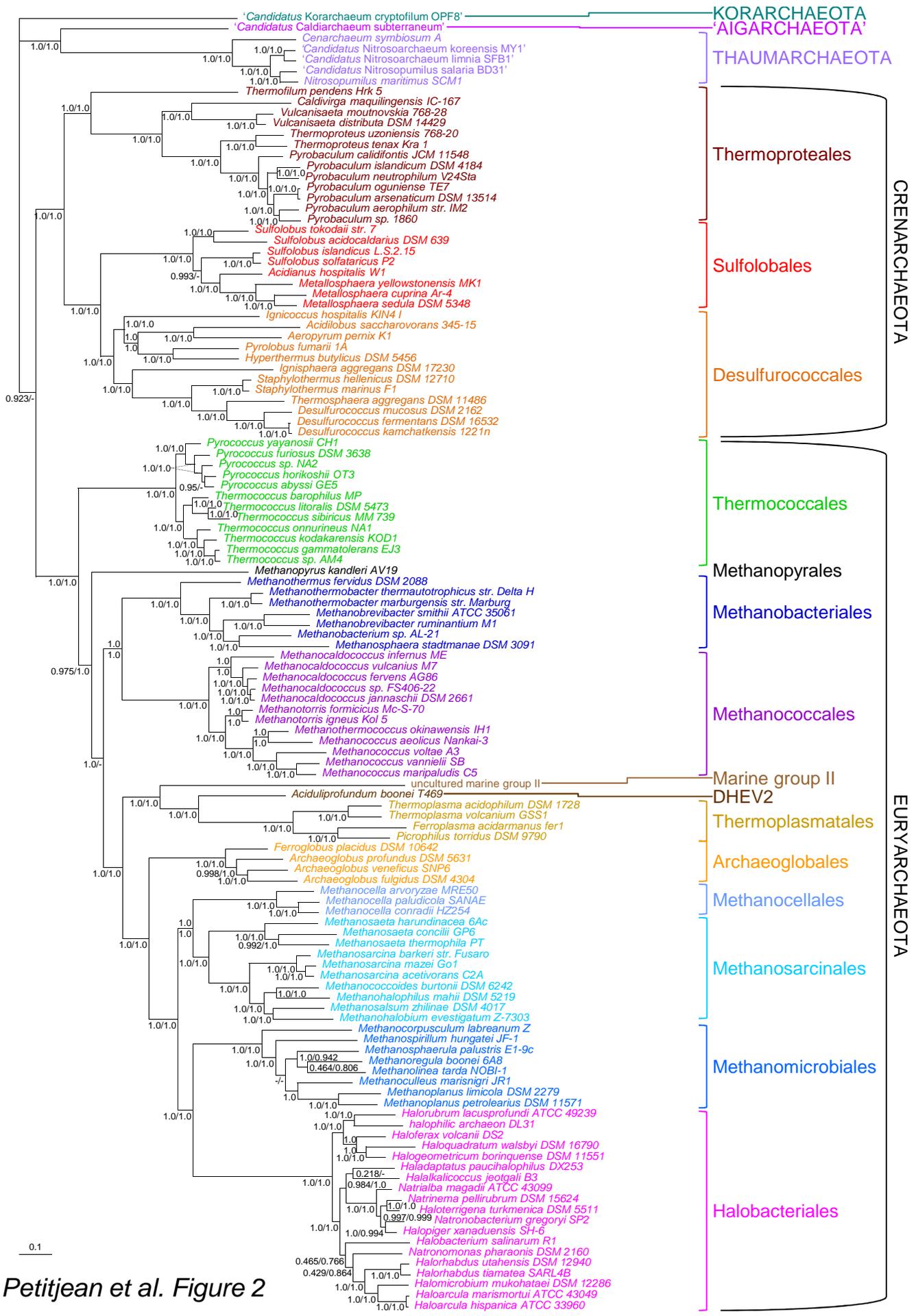
835

836

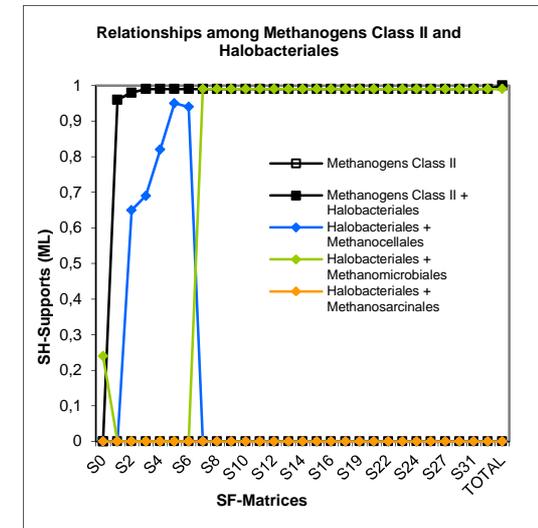
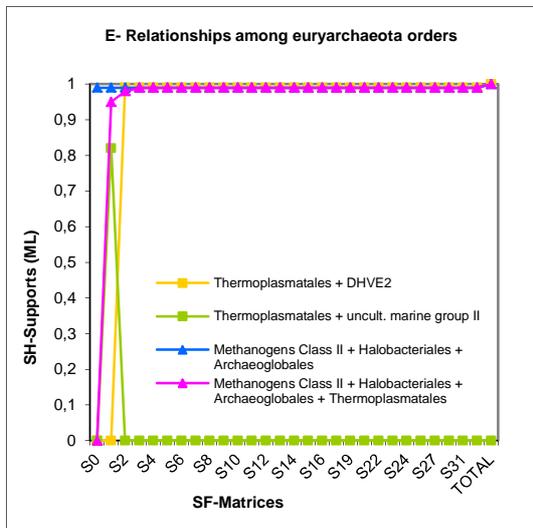
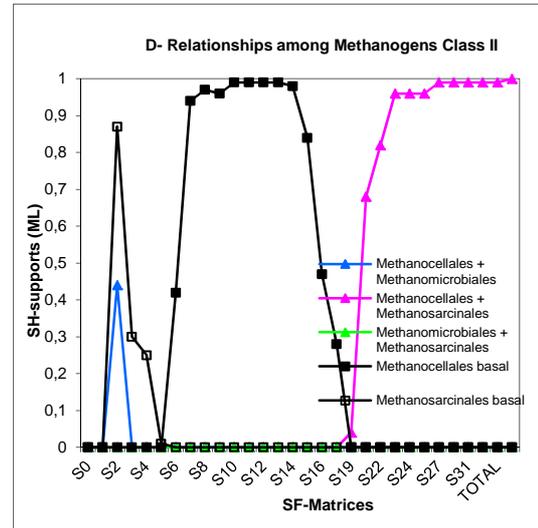
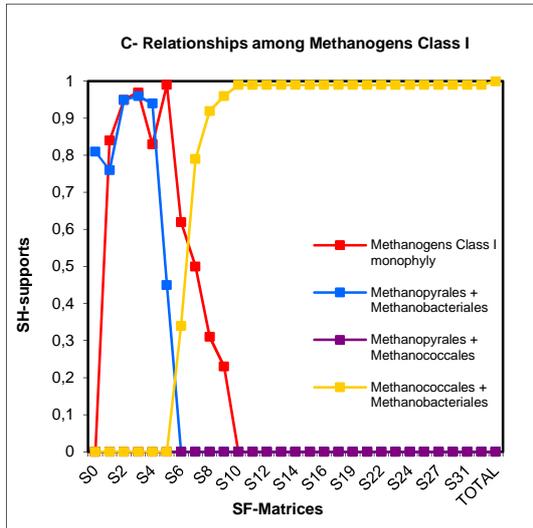
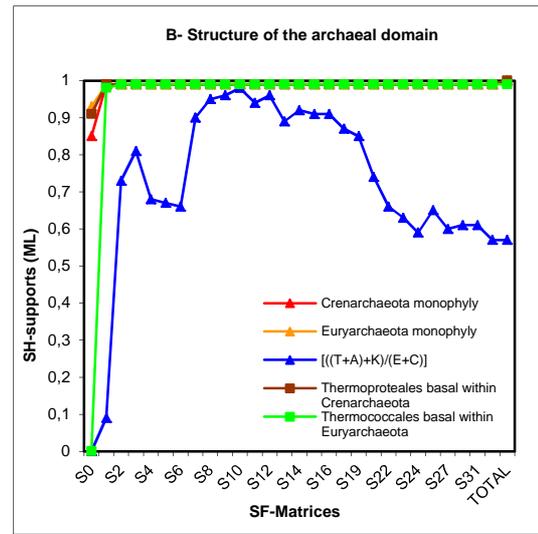
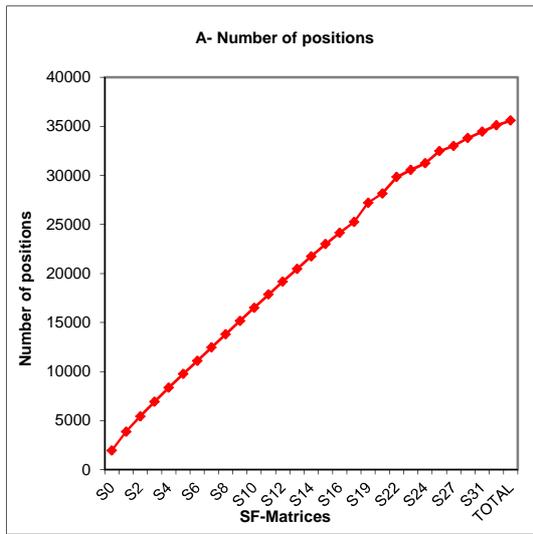


Petitjean et al. Figure 1

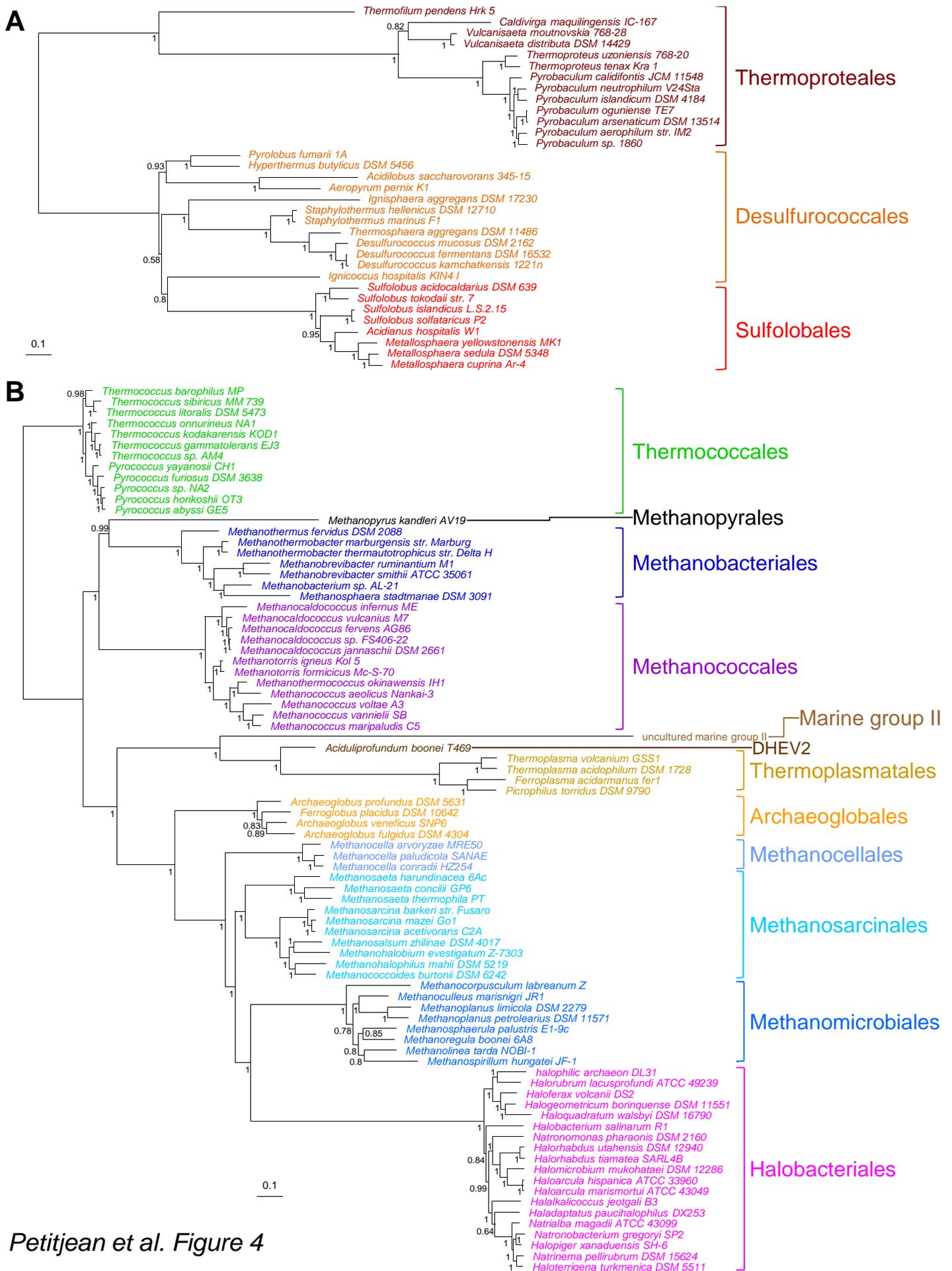




Petitjean et al. Figure 2



Petitjean et al. Figure 3



Petitjean et al. Figure 4

C. Synthèse et éléments de discussion

Les phylogénies obtenues en maximum de vraisemblance à partir des trois supermatrices constituées des protéines ribosomiques, du système de transcription et des 200 nouveaux marqueurs ont des topologies congruentes dans l'ensemble, ce qui est la preuve d'un signal phylogénétique consistant dans ces différents jeux de données. On retrouve la monophylie des différents grands groupes, comme les Euryarchaeota, Crenarchaeota et Thaumarchaeota, comme celle de la plupart des ordres de crenarchées et d'euryarchées. De plus, ces phylogénies montrent des supports aux nœuds assez élevés. Une différence subsiste en ce qui concerne la position de *M. kandleri*, à la base des euryarchées dans l'arbre construit à partir des protéines intervenant dans la transcription, alors qu'il se place au sein des euryarchées dans les autres phylogénies. Cette incongruence est attendue étant données les analyses précédentes qui montraient une attraction de longues branches dans le système de transcription pour cette espèce (Brochier, Forterre, and Gribaldo 2004). Ces résultats sont donc tout à fait en accord avec les résultats des études précédentes, et plus particulièrement avec la phylogénie de Brochier-Armanet et collaborateurs de 2011 (Brochier-Armanet, Forterre, and Gribaldo 2011). Le signal global de ces trois jeux de données n'étant pas contradictoire, nous avons construit deux supermatrices à partir de l'ensemble de ces 273 jeux de données, une avec la totalité des jeux de données et une avec seulement 179 jeux de données contenant moins de 10 espèces manquantes. Les topologies obtenues étaient, comme attendu, en accord avec les résultats précédents. Nous avons alors mené une analyse de ces jeux de données par désaturation par sélection de sites évoluant de plus en plus rapidement pour estimer l'impact de la vitesse d'évolution des sites sur les relations entre archées. Cette analyse a été conduite sur la supermatrice contenant les 179 jeux de données avec peu de données manquantes pour ne pas biaiser la répartition taxonomique dans les différentes supermatrices résultantes et sans les lignées d'archées de tailles nanométrique dont le cas a été traité par d'autres analyses. Les phylogénies obtenues montrent un signal très fort pour la monophylie des Euryarchaeota et des Crenarchaeota, et pour un regroupement des Thaumarchaeota, Aigarchaeota et Korarchaeota, bien que l'ordre de divergence entre ces phyla ne puisse pas être résolu en l'absence de groupe extérieur. Il est aussi intéressant de noter que les premiers ordres divergeant chez les Euryarchaeota et chez les Crenarchaeota sont, respectivement, les *Thermococcales* et les *Thermoproteales*.

Des supermatrices ont ensuite été construites contenant uniquement les Crenarchaeota et les Euryarchaeota afin d'analyser plus précisément les relations au sein de ces deux phyla. Les résultats obtenus correspondent à ce qui était attendu d'après les résultats des analyses sur la totalité des archées. En ce qui concerne les Crenarchaeota, la monophylie des *Thermoproteales* et des

Sulfolobales est retrouvée, mais la monophylie du genre *Sulfolobus* ne l'est pas. Dans ces analyses, la monophylie des *Desulfurococcales* n'est pas retrouvée car *I. hospitalis* (l'hôte de *N. equitans*) se place à la base des *Sulfolobales*. Cette position peut être le signe d'une paraphylie de ce groupe, indiquant que l'ordre des *Desulfurococcales* n'existerait pas, ou d'un artefact dû à une attraction de longues branches ou à des transferts horizontaux non détectés lors des analyses préliminaires. Les résultats préliminaires de la désaturation par sélection de gènes montrent la monophylie des *Desulfurococcales* dans les premières matrices, i.e. rassemblant les gènes évoluant le plus lentement. La paraphylie de ce groupe pourrait donc être plutôt le résultat d'un biais de type attraction de longues branches qu'une réalité biologique. L'espèce *A. saccharovorans*, à partir de laquelle a été proposé récemment l'ordre des *Acidilobales* (Prokofeva et al. 2009), se place en groupe frère d'*A. permix* et au sein des *Desulfurococcales* avec des supports maximaux, ce qui tendrait à prouver définitivement que cette espèce ne représente pas un nouvel ordre de crenarchées. De même que pour les Crenarchaeota, on retrouve la monophylie des ordres d'Euryarchaeota. En ce qui concerne les points non résolus de la phylogénie de ce phylum, les méthanogènes classe I (*Methanobacteriales*, *Methanococcales* et *Methanopyrales*) sont monophylétiques et avec un support maximal, ce qui n'était pas établi jusqu'à présent. De plus, les *Methanobacteriales* et les *Methanopyrales* sont groupes frères avec un très fort support, et ces résultats sont confirmés par les analyses de désaturation par sélection de sites, sur les matrices contenant les sites évoluant le plus lentement. Nous proposons donc une redéfinition taxonomique de ce groupe en lui assignant le nom de « Methanohoma », du grec « ὁμάς » pour « groupe ». On retrouve également la monophylie du groupe Thermoplasmatales – DHVE2 - groupe II d'euryarchées, les deux premiers groupes semblant être plus poches entre eux. Néanmoins, un plus grand nombre de séquences génomiques des groupes DHVE2 et groupe II sera nécessaire pour confirmer ces relations. Un autre point non résolu concerne les relations entre les méthanogènes classe II (*Methanomicrobiales*, *Methanosarcinales* et *Methanocellales*) et les *Halobacteriales*. Nous observons ici que ces quatre ordres forment un groupe monophylétique et bien soutenu dans la plupart des analyses. Par contre, les relations au sein de ce groupe ne sont pas bien résolues. D'après la phylogénie obtenue sans désaturation, les *Halobacteriales* et les *Methanomicrobiales* formeraient un groupe monophylétique, alors que les analyses de désaturation par sélection de sites sur les sites évoluant lentement sont en faveur d'un regroupement entre *Halobacteriales* et *Methanocellales*.

La question de la position phylogénétique des archées de petite taille est un problème majeur. La plupart des phylogénies inférées ces dernières années les placent dans un groupe à la base de l'arbre des archées ou à la base des euryarchées, mais sur de très longues branches. Leur évolution rapide accompagnée d'une réduction de la taille du génome explique facilement pourquoi

une attraction de longues branches peut être la cause de cette position basale et groupée. Néanmoins, cet artefact s'est révélé être très fort et ce groupe reste très difficile à casser pour positionner ces espèces au sein des archées. Des analyses de désaturation par sélection de gènes avec différentes sélections d'espèces (*i.e.* introduction indépendante de chaque espèce à évolution rapide : *N. equitans*, les deux espèces du genre '*Ca. Parvarchaeum*', '*Ca. Micrarchaeum acidiphilum*', et les trois espèces de l'ordre *Nanohaloarchaea*) sont en cours. Cette approche avait déjà montré que *N. equitans* est probablement groupe frère des Thermococcales (Brochier et al. 2005) et nos résultats préliminaires confirment que les *Nanohaloarchaea* sont groupe frère des *Halobacteriales* et que les deux groupes d'ARMAN ('*Ca. Parvarchaeum*' et '*Ca. Micrarchaeum*') pourraient se placer au sein des euryarchées mais séparément. Ces analyses tendent à confirmer que le regroupement de ces sept espèces est artéfactuel dans certaines analyses phylogénétiques, même s'il reste très difficile de les repositionner avec certitude dans la phylogénie des archées. Le fait que ces archées nanométriques aux génomes particuliers ne soient très probablement pas des groupes basaux confirmerait aussi le fait que l'ancêtre des archées n'aurait pas ces caractéristiques génomiques, mais serait au contraire un organisme complexe.

En conclusion, cette analyse montre que le signal phylogénétique porté par les protéines ribosomiques et les sous-unités de l'ARN polymérase n'est pas le résultat d'une évolution particulière de ce système puisqu'on retrouve le même signal avec un grand nombre de marqueurs impliqués dans de nombreux processus cellulaires différents et non liés fonctionnellement. Néanmoins, l'ajout de 200 nouveaux marqueurs permet une amélioration sensible de la reconstruction de la phylogénie des archées en termes de soutien et permet d'apporter un nouvel éclairage sur certains points problématiques. Bien que ces nouveaux marqueurs ne soient pas présents dans la totalité des génomes d'archées, il est possible d'en extraire un signal phylogénétique fiable. Il semble donc qu'un génome « cœur » (« core genome » en anglais) existe et soit hérité majoritairement de manière verticale chez les archées, même s'il peut y avoir des pertes différentielles et HGT ponctuels de certaines protéines dans les différentes lignées. Cette idée avait été proposée par Gribaldo et Brochier-Armanet en 2006 (Simonetta Gribaldo and Brochier-Armanet 2006) et semble confirmée par cette étude. Il est important de souligner que nous avons pris pour point de départ des génomes de thaumarchées et de l'aigarchée, il est donc possible que d'autres marqueurs intéressants soient mis en évidence par l'analyse de génomes d'autres phyla d'archées, et en particulier des lignées non-cultivées à l'heure actuelle et qui constituent une grande part de la diversité du troisième domaine du vivant.

Chapitre 2 : Positionnement de la racine de l'arbre des Archaea.

A. Introduction

La position de la racine de l'arbre des archées est un sujet essentiel non seulement pour la phylogénie des archées mais pour celle du vivant dans son ensemble. En effet, résoudre les relations de parenté entre les phyla d'archées est important pour comprendre l'évolution de ce domaine, mais aussi ses relations avec les deux autres domaines du vivant, bactéries et eucaryotes. En 1990, la première divergence chez les archées a été placée entre Crenarchaeota et Euryarchaeota, définis comme deux phyla par Woese et collaborateurs sur la base des premières phylogénies d'ARNr SSU (C R Woese, Kandler, and Wheelis 1990) (Figure 4). La découverte successive de plusieurs lignées se plaçant près de la base de l'arbre des archées a peu à peu remis en cause cette dichotomie. Ces nouveaux phyla, les Korarchaeota (Barns et al. 1996), les Nanoarchaeota (H. Huber et al. 2002), les Thaumarchaeota (Brochier-Armanet et al. 2008) et les Aigarchaeota (Nunoura et al. 2011), ont multiplié les positions possibles de la racine de l'arbre. Par exemple, les premières phylogénies basées sur l'ARNr SSU et les protéines ribosomiques comprenant la nanoarchée *Nanoarchaeum equitans* plaçaient la racine entre cette espèce et le reste des archées (Waters et al. 2003). De la même manière, les analyses phylogénomiques du génome de la première thaumarchée séquencée, *Cenarchaeum symbiosum*, racinaient l'arbre des archées entre cette espèce et le reste, en utilisant les eucaryotes comme groupe extérieur (Brochier-Armanet et al. 2008). D'autres études phylogénomiques concernant les relations entre les trois domaines du vivant ont été conduites ces dernières années et la racine de l'arbre des archées a été proposée successivement comme étant à la base de chacun des phyla (cf. Introduction B.2.b). Ces analyses ont été conduites avec différentes méthodes et échantillonnages taxonomiques, et avec différents groupes extérieurs (bactéries et/ou eucaryotes), les rendant difficilement comparables. *A contrario*, les marqueurs utilisés sont souvent centrés sur les protéines ribosomiques et, par conséquent, les jeux de données sont assez proches en termes de composition protéique même s'ils ne sont pas analysés de la même façon. Enfin, pour la plupart ces études ont pour but de résoudre l'ensemble des relations entre les trois domaines du vivant, mais ne sont pas centrées sur la position de la racine de l'arbre des archées. D'autre part, une grande partie des phylogénies des archées ont été inférées sans groupe extérieur, ce qui est très important pour pouvoir résoudre la phylogénie interne des archées, mais ne permet pas de raciner l'arbre (Brochier-Armanet, Forterre, and Gribaldo 2011). Nous pensons que les relations anciennes entre les archées doivent être traitées pour elles mêmes, avec un set de marqueurs adapté ainsi qu'un échantillonnage taxonomique représentatif de leur diversité et un choix judicieux du groupe

extérieur.

Nous présentons ici une analyse de la racine de la phylogénie des archées basée sur 70 marqueurs dont 38 ne sont pas des protéines ribosomiques, utilisant des séquences homologues bactériennes comme groupe extérieur. A partir des 200 nouveaux marqueurs définis préalablement pour l'analyse de la phylogénie des archées présentée dans le Chapitre 1, nous avons trouvé 81 protéines ayant des homologues bactériens. Nous avons construit des jeux de données contenant les séquences d'archées et les séquences bactériennes homologues issues de 117 génomes représentatifs de la diversité bactérienne. La phylogénie de chaque protéine a été inférée et analysée afin d'exclure les jeux de données pour lesquels l'homologie des séquences bactériennes était le résultat de transferts horizontaux. 38 jeux de données ont été conservés. Le même travail a été fait à partir des protéines ribosomiques, et 32 ont été conservées. A partir de ces deux groupes de protéines, nous avons tout d'abord construit deux supermatrices différentes, afin d'analyser leur signal indépendamment, puis une troisième supermatrice contenant l'ensemble des 70 jeux de données. Une attention particulière a été donnée aux espèces d'archées de taille nanométrique des genres *Nanoarchaeum*, '*Ca. Micrarchaeum*', '*Ca. Parvarchaeum*' et les *Nanohaloarchaea*, connues pour être sujets aux LBA, ce qui est particulièrement problématique dans une analyse comprenant un groupe extérieur. Enfin, nous avons utilisés deux méthodes de désaturation de notre jeux de données, par gène et par site, afin d'analyser l'impact des différences de vitesse d'évolution sur la position de la racine des archées.

Ce travail est présenté dans l'article « *Phylogenomic Analysis Pinpoints the Root of the Domain Archaea and Supports the Foundation of the New Kingdom Proteoarchaeota* » (Petitjean, Deschamps, Brochier-Armanet, López-García, Moreira, en préparation).

B. Manuscrit de l'article 2 : « Phylogenomic Analysis Pinpoints the Root of the Domain Archaea and Supports the Foundation of the New Kingdom Proteoarchaeota »

Phylogenomic Analysis Pinpoints the Root of the Domain Archaea and Supports the Foundation of the New Kingdom Proteoarchaeota

Céline Petitjean¹, Philippe Deschamps¹, Céline Brochier-Armanet², Purificación López-García¹, and David Moreira^{1,*}

¹Unité d'Ecologie, Systématique et Evolution, CNRS UMR 8079, Université Paris-Sud, 91405 Orsay Cedex, France

²Université de Lyon, Université Lyon 1, CNRS, UMR5558, Laboratoire de Biométrie et Biologie Evolutive, 43 boulevard du 11 novembre 1918, F-69622 Villeurbanne, France. Tel.: +33 4 26 23 44 76; fax: +33 4 72 43 13 88

*Author for Correspondence: David Moreira, Unité d'Ecologie, Systématique et Evolution, CNRS UMR 8079, Université Paris-Sud, 91405 Orsay Cedex, France, +33 1 69157608, +33 1 69154697, david.moreira@u-psud.fr

Data deposition: All sequence alignments used in this work are available upon request to the corresponding author.

Abstract

The first phylogenetic trees of the domain Archaea, based on the analysis of 16S rRNA sequences, showed a deep division between two major groups, which were classified as the kingdoms Euryarchaeota and Crenarchaeota. This bipartite view has been challenged in recent years thanks to the discovery of several deeply-branching new archaeal lineages, in particular the Thaumarchaeota, Aigarchaeota, Nanoarchaeota, and Korarchaeota, which have also been given the same taxonomic status of kingdoms. However, the phylogenetic position of some of these lineages has been controversial. In this sense, one problem has been that very often the phylogenetic analyses of the Archaea has been carried out without outgroup, making very difficult to determine if these taxa actually represent independent archaeal lineages of the same level as the Euryarchaeota and Crenarchaeota. We have addressed this question by reconstructing archaeal phylogenetic trees rooted on bacterial sequences. These trees were based on conserved protein markers commonly used (32 ribosomal proteins) and on 38 new markers identified through phylogenomic analysis. We thus gathered a total of 70 markers that we analyzed as a concatenated dataset. Our trees consistently placed the root of the archaeal tree between the Euryarchaeota (including the Nanoarchaeota and other fast-evolving lineages) and the rest of archaeal species. We thus propose the erection of a new kingdom, the Proteoarchaeota, to contain these lineages. This has to be accompanied by the relegation to a lower taxonomic rank (class or order) of the lineages previously classified as kingdoms (Crenarchaeota, Thaumarchaeota, Aigarchaeota, and Korarchaeota) and now encompassed within the Proteoarchaeota.

Key words: Archaea, Euryarchaeota, Proteoarchaeota, root, phylogenomics.

Introduction

Despite the fact that several archaeal species had been isolated as soon as in the 30's (Barker, 1936), the recognition of their nature as members of an independent domain of life had to wait for 40 years, when Woese and Fox published the proposal for the division of life into three primary kingdoms: Bacteria, Archaeobacteria and Eukaryotes, based on the analysis of small subunit rRNA sequences (SSU rRNA) (Woese and Fox 1977). Later on, to emphasize on the difference between the two major groups of prokaryotes (Bacteria and Archaeobacteria), these three kingdoms were reclassified as the domains Bacteria, Archaea and Eucarya (Woese et al., 1990). Within the Archaea, the first SSU rRNA phylogenies supported the separation of two groups, one containing methanogenic species and another formed by thermoacidophilic ones (Fox et al., 1980). This was confirmed by subsequent analyses with a richer taxonomic sampling, leading to the division of the Archaea into two kingdoms: the Crenarchaeota, all hyperthermophilic, and the Euryarchaeota, containing species with a variety of phenotypes (hyperthermophilic, mesophilic, methanogenic, halophilic, etc.) (Woese et al., 1990). However, a number of discoveries have challenged this simple bipartite view in recent years.

First, environmental SSU rRNA sequence analyses revealed the existence of archaeal species related to the crenarchaeota thriving in mesophilic environments such as the open ocean (DeLong, 1992; Fuhrman et al., 1992), soils (Jurgens et al., 1997), lakes (Schleper et al., 1997) and both cold and hot terrestrial springs (Barns et al., 1996). Phylogenetic analysis of conserved genes involved in translation and some differences in gene content with the classical hyperthermophilic Crenarchaeota led to propose that these archaeal species may define a new archaeal phylum, the Thaumarchaeota, independent from the two phyla recognized until now, the Crenarchaeota and the Euryarchaeota (Brochier-Armanet et al., 2008). More recently, a distant relative of the Thaumarchaeota, the species *Candidatus 'Caldiarchaeum subterraneum'*, was proposed to define an additional new phylum - Aigarchaeota- based on its distinct gene content (Nunoura et al., 2011). Similar arguments were used to suggest that the species *Candidatus 'Korarchaeum cryptofilum'* represented the first member of a new phylum Korarchaeota, very distantly related to the Crenarchaeota (Barns et al., 1996; Elkins et al., 2008). An even more divergent case was found with the discovery of the hyperthermophilic *Nanoarchaeum equitans*, a tiny symbiont of the crenarchaeote *Inginicoccus hospitalis*. The phylogenetic analysis of concatenated ribosomal proteins suggested that this species diverged before the separation of the Crenarchaeota and the Euryarchaeota, thus defining a very ancient new archaeal phylum -the Nanoarchaeota- (Waters et al., 2003). Although being classified within the Euryarchaeota, a very deep branching has also been suggested for other ultrasmall archaea of the genera

Candidatus 'Micrarchaeum' and 'Parvarchaeum', using as argument their unusual gene content (Baker et al., 2010).

Thus, the last decade has seen a multiplication of deep-branching groups within the Archaea and, as a consequence, of possibilities to place the root of the archaeal tree, namely, the first divergence within this domain of Life. The first suggestion placed the root between the Crenarchaeota and the Euryarchaeota (Woese et al., 1990). However, the subsequent inclusion of the new deep-branching groups, in particular the Nanoarchaeota, has challenged this view. For example, the first phylogenetic trees incorporating this lineage, based on rRNA or ribosomal protein sequences, supported a rooting between *N. equitans* and the rest of archaea (Waters et al., 2003). However, other analyses rooted the archaeal tree on the branch leading to the Thaumarchaeota (Brochier et al., 2008). More recently, phylogenetic analyses using heterogeneous sequence evolution models aimed at determining the precise archaeal origin of the eukaryotic nucleocytoplasm have placed the root of the archaeal domain on the *N. equitans* branch or within the Euryarchaeota (Cox et al., 2008). These different analyses applied different tree reconstruction methods, different models of sequence evolution, different archaeal sequence samplings and different outgroup sequences (eukaryotic and/or bacterial), making them difficult to be compared. On the other hand, most archaeal phylogenetic analyses with a wide taxonomic archaeal sampling did not include outgroup sequences, so they only produced unrooted phylogenies (e.g., Brochier-Armanet et al., 2011). Therefore, there is a need for phylogenetic analyses specifically aimed at rooting the archaeal tree. For that, a wide representation of all archaeal phyla as well as a good sampling of outgroup sequences are necessary. In this work, we have used this approach to reconstruct archaeal phylogenetic trees rooted with a bacterial outgroup. These trees were based on the classical ribosomal protein dataset (32 proteins) and on a collection of 38 new conserved proteins identified by phylogenomic analysis. Our phylogenetic analyses supported the rooting between the Euryarchaeota (including the Nanoarchaeota and the genera *Micrarchaeum* and *Parvarchaeum*) and the rest of archaea (Crenarchaeota, Thaumarchaeota, Aigarchaeota and Korarchaeota). This deep division into two major groups prompted us to propose a re-classification of several major archaeal lineages from their current status of phyla or divisions into the rank of classes. This would make the whole archaeal taxonomy much more homogeneous.

Materials and Methods

Detection of Proteins Widespread in Archaea and Bacteria

We gathered the protein sequences coded by all publicly available complete genome sequences of archaeal species as well as those of 117 bacterial species into a local

database (for the list of organisms, see supplementary table S1, Supplementary Material online). We then carried out sequence similarity searches against the archaeal sequences in this database using BLAST (Altschul 1997) and all the protein sequences coded in the genomes of *Nitrosopumilus maritimus*, *Cenarchaeum symbiosum* and *Caldiarchoaeum subterraneum* as queries. The BLAST results were filtered to retain only those containing representatives at least 4 archaeal sequences. The 4174 proteins thus identified were submitted to preliminary phylogenetic analysis (see below). The phylogenetic trees obtained were checked manually to retain only those where the different archaeal phyla were monophyletic and without evidence for horizontal gene transfer (HGT) events.

200 datasets were retained, for which we incorporated the respective homologous bacterial sequences and reconstructed new phylogenetic trees, which we used to finally retain 38 proteins for which the trees did not show evidence for HGT between archaea and bacteria (see the list in supplementary table S2, Supplementary Material online). To decrease the calculation time of the subsequent phylogenetic analyses, we selected 27 taxonomically diverse bacteria to be used as outgroup.

Phylogenetic Analyses

Sets of homologous protein sequences were aligned using MAFFT (Kato 2013). Conserved positions in the alignments were detected using BMGE with default parameters (Crisuolo 2010) and verified by hand using the program NET of the MUST package (Philippe 1993). Maximum likelihood phylogenetic trees were reconstructed upon each individual protein or different concatenated datasets with RaxML v.7.2.4 (Stamatakis 2006) and the PROTGAMMALGF model. Tree robustness was estimated using the Rapid Bootstrapping method implemented in RaxML. Bayesian inference analyses were carried out using MrBayes 3.2.1 (Ronquist et al., 2012) with a mixed model of amino acid sequence evolution and a Gamma distribution with six discrete categories to accommodate for among site rate variation.

Mutational Desaturation Analyses

We applied a site-by-site mutational desaturation approach using SlowFaster (Kostka 2008) to construct 10 multiple sequence alignments containing sites with increasing evolutionary rates. We also carried out a marker-by-marker desaturation analysis based on the average evolutionary rates of the different markers analyzed. For that, we selected 22 pairs of species (10 archaea-archaea, 9 archaea-bacteria, and 3 bacteria-bacteria, see supplementary table S3, Supplementary Material online) and we computed the species-to-species distance for each pair based on the ML phylogenetic tree for each individual marker. We then obtained the average value of those distances, which we attributed as the evolutionary rate of the

marker. As for the previous approach, we constructed 10 concatenations, beginning with one containing the slowest evolving markers and progressively adding more and more fast evolving ones. All the different concatenations were analyzed by maximum likelihood using with RaxML v.7.2.4 (Stamatakis 2006).

Results

Identification of New Markers Widely Distributed in Archaea and Bacteria

It has been proposed that the gene content in thaumarchaeotal genomes is to some extent intermediate between the euryarchaeotal and crenarchaeotal ones (Brochier et al., 2008, Spang et al., 2010, Brochier et al., 2012). By this reason we have used two thaumarchaeotal genomes (*Cenarchaeum symbiosum* and *Nitrosopumilus marinus*) and one from the closely related phylum Aigarchaeota (*Caldiarchaeum subterraneum*) as queries to look for conserved genes widely distributed in Archaea. Using BLAST searches and phylogenetic analysis, we detected 200 proteins with orthologs in most of the available archaeal complete genome sequences. Among them, 68 had homologues in bacterial species. We reconstructed maximum likelihood (ML) phylogenetic trees for all these proteins to verify that their presence in both Archaea and Bacteria was not due to horizontal gene transfer (HGT) events between these two domains. We also checked that these proteins were present in at least 26 bacterial phyla to have a sufficiently diverse set of outgroup sequences. This yielded 38 well conserved proteins widely distributed in Archaea and Bacteria (supplementary table S1, Supplementary Material online). We excluded eukaryotes from all our analyses because of the possibility that they might have inherited part of their genes from Archaea during their very early evolution (Cox et al., 2008; Embley and Martin 2006; Gribaldo et al., 2010; López-García and Moreira 1999). If this is actually the case, their use as an outgroup would lead to infer an erroneous rooting of the archaeal tree.

In contrast with the limited range of functions of the markers used until now to reconstruct rooted phylogenies of the Archaea (mostly ribosomal proteins), our newly detected proteins are involved in a variety of cellular processes, including the metabolism of amino acids, nucleotides and co-enzymes, post-translational protein modification and other functions (supplementary table S1, Supplementary Material online). This functional diversity is important because it can contribute to retrieve more homogeneous branch lengths among species because it is unlikely that this entire set of very different cellular processes might have an evolutionary rate in some lineages significantly more accelerated than in the rest. Such very unequal evolutionary rates among lineages, which represent a major problem for

phylogenetic reconstruction, are more likely to occur for small sets of proteins involved in a same process.

In addition to these new protein markers, we also included the classical set of translation proteins (32 ribosomal proteins) to construct concatenated sequence datasets for the subsequent phylogenetic analyses.

Rooted Phylogenetic Analyses of Translation Proteins and the Newly Detected Markers

Multi-marker phylogenetic analyses have revealed that several archaeal lineages have very high evolutionary rates, producing very long branches in phylogenetic trees which can lead to long branch attraction artefacts (LBA). This is notably the case for the ultrasmall archaea of the genera *Nanoarchaeum*, *Micrarchaeum*, *Parvarchaeum* and the Nanohaloarchaea (e.g., Brochier et al., 2011). Therefore, we carried out a first set of phylogenetic analyses excluding all these fast-evolving organisms. The first dataset that we analyzed contained the classical set of ribosomal proteins that have traditionally been used to reconstruct rooted and unrooted archaeal phylogenies (e.g., Cox et al., 2008, Elkins et al., 2008, Brochier et al., 2011, Guy and Ettema 2011). We reconstructed an updated phylogenetic analysis of those markers (32 markers and 2560 amino acid sites) incorporating representative species from all major archaeal lineages (with the exception of the fast-evolving ones, see above) to obtain a reference phylogeny. The resulting phylogeny was rooted within the Euryarchaeota (fig. 1). In fact, different euryarchaeotal lineages emerged as paraphyletic branches at the base of the tree, with strong statistical support in Bayesian trees (posterior probabilities PP=0.97-1) but moderate bootstrap support in maximum likelihood trees (BP<70%). This result was in agreement with a similar tree published by Cox et al., (2008). The apical part of our tree was occupied by a strongly supported group (PP=1, BP=100%) containing the Thaumarchaeota, Aigarchaeota, Crenarchaeota and Korarchaeota (a group tentatively defined as the "TACK" superphylum by Guy and Ettema (2011)).

We then carried out a phylogenetic analysis on a concatenation of our new markers (38 markers and 6890 amino acid sites). The Bayesian phylogenetic tree based on this new dataset (fig. 2) was rooted between the Euryarchaeota and the rest of archaeal species, namely the "TACK" group, in sharp contrast with the previous tree based on translation-related proteins. This rooting was robustly supported (PP=1, BP=99-100%). Moreover, the monophyly of most of the archaeal classes was also well supported (PP=1, BP>90%).

Finally, we analyzed a combined dataset contained all the previous markers (translation-related proteins plus our new markers, for a total of 70 markers and 9450 amino acid sites). The resulting tree (fig. 3) was very similar to the one obtained from the analysis of the concatenation of the new markers. In fact, the root was placed between the Euryarchaeota

and the rest of archaeal species with strong support, especially by the Bayesian analysis (PP=1 for all deep nodes in the tree) and moderately by the ML analysis (BP=79% for the monophyly of the Euryarchaeota and 100% for the "TACK" group). The internal phylogeny of the different archaeal lineages was also well supported in the Bayesian tree, even for several difficult-to-resolve relationships. For example, we retrieved full support (PP=1, BP=100%) for the monophyly of the so-called "Methanogen Class I" containing the Methanopyrales, Methanobacteriales and Methanococcales, a relationship that was not supported by the ribosomal proteins, which placed the Methanococcales as sister group of the Thermococcales (fig. 1). Another highly supported (PP=1, BP=99%) interesting result was the position of the Halobacteria as a derived group within the class Methanomicrobia (also known as "Methanogen Class II").

The Position of Long-Branching Archaea

As mentioned above, we excluded from our initial phylogenetic analyses a number of taxa characterized by their very long branches. These included *Nanoarchaeum*, *Micrarchaeum*, *Parvarchaeum* and the Nanohaloarchaea. These species exhibited a very large amount of missing data in our concatenated alignments (>50%), in agreement with their very small genome sizes (especially for *N. equitans*, with only 540 protein-coding genes, and *Micrarchaeum* and *Parvarchaeum*, with ~1000 protein-coding genes). This, in addition to their long branches that suggest an accelerated evolutionary rate, made them very prone to potential LBA artefacts (Roure et al., 2013). Moreover, LBA can be exacerbated by the inclusion of distant outgroup sequences, leading to a basal emergence of the long-branching taxa attracted by the outgroup (Philippe and Laurent 1998). In fact, although poorly supported, the ML tree based on our complete set of markers and rooted using the bacterial sequences as outgroup showed a basal emergence of these long-branching taxa within a monophyletic group (BP=62%, supplementary fig. S1, Supplementary Material online). This would support that the root of the archaeal tree is located between these species and the rest of archaea. However, previous analyses have provided strong support for alternative placements of these species, in particular for *Nanoarchaeum* as sister group of the Thermococcales (Brochier et al., 2005) and the Nanohaloarchaea as sister group of the Halobacteria (Narasingarao et al., 2012). We thus tested different ways to reduce the potential LBA responsible of the basal emergence of the long-branching archaea. First, we carried out an unrooted phylogenetic analysis excluding the bacterial sequences. In the resulting ML tree all the ultrasmall archaea branched again as a strongly supported monophyletic group placed between the Crenarchaeota and the Euryarchaeota (supplementary fig. S2, Supplementary Material online). This suggested that the removal of the distant bacterial outgroup sequences was not enough to alleviate the possible LBA

artefact responsible of their basal emergence. We thus re-analyzed the same dataset using a mixed model of amino acid sequence evolution and Bayesian inference, which is known to be very often more robust against LBA than ML methods (Philippe et al., 2011). Indeed, the Bayesian tree showed a very different topology, with the four groups on long branches well nested at different locations within the Euryarchaeota (fig. 4). Nanoarchaeota branched as sister group of the Thermococcales, the genera *Micrarchaeum* and *Parvarchaeum* as sisters of the Thermoplasmata, and the Nanohaloarchaea as sisters of the Halobacteria. These results supported that these taxa can be considered as *bona fide* euryarchaeotal species and that, therefore, the root of the archaeal tree does not lie between them and the other archaea.

Desaturation Analyses

Since the root corresponds to the most ancient node in the archaeal tree, it can be one of the most affected by mutational saturation, namely the loss of phylogenetic signal due to the accumulation of successive mutations on variable sequence sites. To test if this was the case, we have analyzed our global dataset (translation related proteins plus our new markers) using two desaturation strategies based on site selection and marker selection, respectively. For the first one, we calculated the evolutionary rate at each position of the global alignment and constructed 10 concatenations of increasing size. The first contained only the slowest evolving sites and the following a progressively increasing number of more and more fast evolving sites. The ML phylogenetic tree based on the first dataset (with only 939 sites) placed the root within the Euryarchaeota, between the Thermococcales and the rest of archaeal species. However, all the other datasets recovered the same rooting pattern as the complete dataset, namely, the separation between the Euryarchaeota and the other archaea (see supplementary figs. S3-S12, Supplementary Material online). Whereas the monophyly of the "TACK" group was strongly retrieved by all datasets (BP=100% except for the first dataset, with BP=83%), the monophyly of the Euryarchaeota was less supported, especially by the small- and middle-size datasets. This probably reflected the influence of the ribosomal proteins, which tended to support a root within the Euryarchaeota (fig. 1).

The second desaturation strategy took into account the evolutionary rate of each marker to construct a series of 10 concatenations from the slowest evolving to the fastest evolving ones (see Materials and Methods). As in the site-by-site desaturation analysis, all datasets retrieved the monophyly of the "TACK" group with maximal support (BP=100%, supplementary fig. S13-S22, Supplementary Material online). However, the position of the root was unstable, with some of the smallest datasets supporting a root within the Euryarchaeota. This probably reflected the predominance of ribosomal proteins in these datasets. Interestingly, the last five datasets retrieved the root between the Euryarchaeota

and the other archaea, with highest support (BP=94-97%) in the intermediate datasets (49 and 56 proteins). They probably corresponded to the datasets with the optimum number of informative sites before an increase of noise provided by the fastest-evolving proteins, which led to a decrease of the statistical support in the final datasets.

In summary, the two desaturation strategies, site-by-site and gene-by-gene, produced similar results, with most trees supporting the rooting between the Euryarchaeota and the "TACK" group. The global stability of this result suggested that it is not due to any particular group of genes or site category.

Discussion

The Root of the Archaeal Tree

Since the initial proposal by Woese and co-workers that the deepest division in the archaeal domains was between the two kingdoms Crenarchaeota and Euryarchaeota (Woese 1990), many other alternatives have been advanced, most of them linked to the discovery of new deeply-branching archaeal lineages. As summarized in the Introduction, this has been the case especially for the hyperthermophilic *Nanoarchaeum equitans* (Waters et al., 2003, Cox et al., 2008), the Thaumarchaeota (Brochier et al., 2008), and the acidophilic genera *Candidatus* 'Micrarchaeum' and 'Parvarchaeum' (Baker et al., 2010). Di Giulio has repeatedly argued for the rooting on the Nanoarchaeota, to even propose that this group should be considered as a "living fossil" (Di Giulio 2006). This was based not only on phylogenetic analyses of rRNA sequences (Branciamore 2008) but also on the unusual discovery of some genes split in two fragments in the genome of *N. equitans* (Randau 2005), which was interpreted as an ancestral character according to the "introns early" hypothesis (Di Giulio 2008). However, several findings have undermined this proposal. Split genes have also been found in another small archaeon, *Micrarchaeum acidiphilum*, which is phylogenetically unrelated to *N. equitans* (Baker et al., 2010). In addition, phylogenetic analyses based on individual conserved proteins and on concatenated translation-related proteins strongly support that Nanoarchaeota are not basal but sisters to the Thermococcales (Brochier 2005; Brochier 2011). Likewise, a very deep-branching position for the genera *Micrarchaeum* and *Parvarchaeum* based on their gene content with both typical euryarchaeotal and crenarchaeotal genes (Baker et al., 2010) has not been validated by phylogenetic analysis of translation-related proteins, which placed them nested within the Euryarchaeota (Brochier 2011). Finally, the third group of long-branching archaea, the halophilic Nanohaloarchaea, has been shown to robustly branch close to the Halobacteria (Narasimgarao et al., 2012). Our analyses with the complete set of markers are in agreement with these results. The unrooted Bayesian tree showed all the long-branching archaea well nested within the Euryarchaeota,

with the Thermococcales + Nanoarchaeota as the first group to diverge (fig. 4). Interestingly, when bacterial sequences were added to root this tree, all those long-branching archaea emerged together at the base of the archaeal tree, strongly suggesting a LBA artefact.

Thus, the atypical characteristics found in several fast-evolving ultrasmall archaea are most likely derived characters rather than ancestral features. One example concerns the split genes of *Nanoarchaeum* and *Micrarchaeum*, which have likely evolved by convergence in these two lineages. Massive genome size reduction has also occurred in parallel in these organisms as well as in *Parvarchaeum*, and, to a lesser extent, in the Nanohaloarchaea, all of them having gene numbers smaller than those of their respective closest relatives. This has been accompanied by a large acceleration of evolutionary rate, as attested by the long branches exhibited by all these species. It may be possible that anomalous characters, such as the split genes, are side products of that evolutionary acceleration and genome reduction. In fact, convergent acquisition of exceptional features has already been noticed in fast-evolving highly reduced genomes. One outstanding example is the migration of the rRNA genes to subtelomeric regions in the chromosomes of the highly reduced nucleomorph genomes of cryptophytes and chlorarachniophytes (Moore and Archibald 2009).

A Bipartite Division of the Domain Archaea and Proposal for the New Archaeal Kingdom Proteoarchaeota

Originally, the domain Archaea was divided into two major groups: the Euryarchaeota and the Crenarchaeota, based on the analysis of 16S rRNA sequences (Woese et al., 1990). These two groups were given the taxonomic rank of Kingdom as the one immediately below that of Domain and to insist on their clear phylogenetic distinctiveness (Woese et al., 1990). However, since the publication of this proposal, several new archaeal lineages have been discovered and some of them have been named using the kingdom suffix 'archaeota' to equal their ranks to that of the Euryarchaeota and Crenarchaeota. This is the case of the Nanoarchaeota (Waters et al., 2003), Thaumarchaeota (Brochier-Armanet et al., 2008), Aigarchaeota (Nunoura et al., 2011), and Korarchaeota (Elkins et al., 2008). The lack of clear criteria to establish a rank for those different lineages has fostered discussion on different aspects, such as the weight that molecular phylogeny has to have on the definition of new major taxa (Gribaldo and Brochier 2012; Garrity and Oren 2012).

Our phylogenetic analyses excluding the long-branching taxa (*Nanoarchaeum*, *Micrarchaeum*, *Parvarchaeum* and the Nanohaloarchaea) strongly supported that the root of the archaeal tree lies between the Euryarchaeota and the rest of archaeal lineages. This, together with the observation that the long-branching archaea can be considered as *bona fide* Euryarchaeota since they branch within this archaeal group when LBA problems are minimized, advocates for a major division of the domain Archaea into two major groups

between which the root is located: the Euryarchaeota and the so-called TACK supergroup. In addition to these phylogenetic considerations, these two groups have similar levels of ecological and evolutionary diversity. As their name evokes (*eurus* meaning wide), Euryarchaeota have for long been known to exhibit a variety of metabolic capacities and to occupy a broad range of habitats (Woese et al., 1990). This is also the case for the TACK supergroup (Guy and Ettema 2011), thriving in high-temperature environments (Crenarchaeota, Aigarchaeota and Korarchaeota) but also in mesophilic ones (Thaumarchaeota) and relying on a large variety of metabolisms. On the other hand, the evolutionary divergence among distant TACK members is very similar to that among distant Euryarchaeota. For example, the average 16S rRNA sequence identity between Thaumarchaeota and Desulfurococcales (Crenarchaeota) is of ~75%, identical to that between Thermococcales and Halobacteria within the Euryarchaeota. The phylogenetic depth, in terms of sequence divergence, for our set of protein markers was also similar for the Euryarchaeota and the TACK group (~65% amino acid sequence similarity for comparisons of the same taxa as above).

The separation of the Euryarchaeota and the TACK group represents the primary split among the known archaeal species and these two groups have comparable ecological and phylogenetic diversities. Thus, it would be logical to give the TACK group the same taxonomic level as the Euryarchaeota, namely a kingdom rank. This would require providing a formal name to the TACK group. We propose to call this new kingdom Proteoarchaeota, making reference to the Greek god of the sea Proteus, able to display many different forms. The same prefix Proteo- was used in the name Proteobacteria also to point to the vast diversity of this bacterial group (Stackebrandt et al., 1988). The erection of the kingdom Proteoarchaeota would entail the relegation in rank of several archaeal lineages that were given a kingdom (or superphylum) level. As mentioned above, this concerns the Nanoarchaeota, Thaumarchaeota, Aigarchaeota and Korarchaeota. They should be reclassified as classes as occurs for the different lineages that compose the Euryarchaeota. Therefore, we propose to apply them the new names Nanoarchaea, Thaumarchaea, Aigarchaea and Korarchaea, with their respective orders Nanoarchaeales, Thaumarchaeales, Aigarchaeales and Korarchaeales (table 1). Likewise, the former kingdom Crenarchaeota should be renamed as a class (Crenarchaea, which would be synonym of Thermoprotei). We consider that this amended scheme is the one that requires the smallest number of taxonomic changes and, at the same time, would be much more homogeneous for the whole archaeal domain than the current mix of different ranks (kingdoms, classes and orders) to refer to lineages that, actually, can be joined into only two major groups. Incidentally, our results reinforce that the Halobacteria branch with high support within the class Methanomicrobia, which turns out to be paraphyletic (fig. 3). Therefore, in order to

avoid the existence of a paraphyletic class, the Halobacteria should be retrograded to the status of order within the class Methanomicrobia. This is in agreement with the evidence supporting that the ancestor of the Halobacteria was a methanogen species that secondarily adapted to hypersaline environments (Nelson-Sathi et al., 2012).

Conclusions

In addition to the classical markers commonly used until now, most of them ribosomal proteins, our phylogenomic survey has allowed identifying 38 additional conserved proteins that can be used to reconstruct phylogenetic trees of the archaea rooted on bacterial homologues. The phylogenetic analyses of the complete set of all those markers (32 ribosomal and 38 new ones) converge to support a deep division of the domain Archaea in two major lineages. One corresponds to the kingdom Euryarchaeota, already defined 23 years ago (Woese et al., 1990) and the second to a miscellaneous collection of lineages that has been tentatively grouped under the informal denomination of TACK supergroup (Guy and Ettema 2011). This second group is clearly monophyletic in our analyses and has a level of phylogenetic diversity comparable to the one of the Euryarchaeota. In addition, the lineages composing it show a large panel of ecological adaptations. These points are not reflected by the taxonomic classification of the Archaea currently used, which gives similar ranks to groups as different in phylogenetic diversity and depth as the extremely wide Euryarchaeota and the much more reduced Nanoarchaeota, Korarchaeota or Aigarchaeota. We think that the best alternative would be to give the two major archaeal lineages the same taxonomic rank. For that, the most parsimonious solution is keeping the kingdom rank already given to the Euryarchaeota and to erect a new kingdom to contain the TACK lineages. We propose to call these new kingdom Proteoarchaeota to highlight its high ecological and phylogenetic diversity.

Supplementary Material

Supplementary figures S1–S22 and tables S1-S3 are available at <http://www.xxx.xxx/>.

Acknowledgments

This work was supported by the French National Agency for Research (EVOLDEEP project, contract number ANR-08-GENM-024-002) and by the Investissement d'Avenir grant (ANR-10-BINF-01-01). C.B-A is member of the Institut Universitaire de France.

Literature Cited

Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.* 1997; 25:3389-402.

Barker, H. A. Studies upon the methane-producing bacteria. *Arch. Mikrobiol.* 1936; 7:420–438.

Baker BJ, Comolli LR, Dick GJ, Hauser LJ, Hyatt D, Dill BD, Land ML, Verberkmoes NC, Hettich RL, Banfield JF. Enigmatic, ultrasmall, uncultivated Archaea. *Proc Natl Acad Sci U S A.* 2010; 107:8806-8811.

Barns SM, Delwiche CF, Palmer JD, Pace NR. Perspectives on archaeal diversity, thermophily and monophyly from environmental rRNA sequences. *Proc Natl Acad Sci USA.* 1996; 93:9188–9193.

Branciamore S, Gallori E, Di Giulio M. The basal phylogenetic position of *Nanoarchaeum equitans* (Nanoarchaeota). *Front Biosci.* 2008; 13:6886-6892.

Brochier C, Gribaldo S, Zivanovic Y, Confalonieri F, Forterre P. Nanoarchaea: representatives of a novel archaeal phylum or a fast-evolving euryarchaeal lineage related to Thermococcales? *Genome Biol.* 2005; 6:R42.

Brochier-Armanet C, Boussau B, Gribaldo S, Forterre P. Mesophilic Crenarchaeota: proposal for a third archaeal phylum, the Thaumarchaeota. *Nat Rev Microbiol.* 2008; 6:245–252.

Brochier-Armanet C, Forterre P, Gribaldo S. Phylogeny and evolution of the Archaea: one hundred genomes later. *Curr Opin Microbiol.* 2011; 14:274-281.

Brochier-Armanet C, Gribaldo S, Forterre P. Spotlight on the Thaumarchaeota. *ISME J.* 2012; 6:227-230.

Cox CJ, Foster PG, Hirt RP, Harris SR, Embley TM. The archaeobacterial origin of eukaryotes. *Proc Natl Acad Sci U S A.* 2008; 105:20356-20361.

Criscuolo A, Gribaldo S. BMGE (Block Mapping and Gathering with Entropy): a new software for selection of phylogenetic informative regions from multiple sequence alignments. *BMC Evol Biol.* 2010; 10:210.

DeLong EF. Archaea in coastal marine environments. *Proc Natl Acad Sci USA.* 1992; 89:5685–5689.

Di Giulio M. *Nanoarchaeum equitans* is a living fossil. *J Theor Biol.* 2006; 242:257-260.

Di Giulio M. The split genes of *Nanoarchaeum equitans* are an ancestral character. *Gene.* 2008; 421:20-26.

Elkins JG, Podar M, Graham DE, Makarova KS, Wolf Y, Randau L, Hedlund BP, Brochier-Armanet C, Kunin V, Anderson I, Lapidus A, Goltsman E, Barry K, Koonin EV, Hugenholtz P, Kyrpides N, Wanner G, Richardson P, Keller M, Stetter KO. A korarchaeal genome reveals insights into the evolution of the Archaea. *Proc Natl Acad Sci U S A*. 2008; 105:8102-8107.

Embley TM, Martin W. Eukaryotic evolution, changes and challenges. *Nature*. 2006; 440:623-630.

Fuhrman JA, McCallum K, Davis AA. Novel major archaeobacterial group from marine plankton. *Nature*. 1992; 356:148–149.

Garrity GM, Oren A. Response to Gribaldo and Brochier-Armanet: time for order in microbial systematics. *Trends Microbiol*. 2012; 20:353-354.

Gribaldo S, Brochier-Armanet C. Time for order in microbial systematics. *Trends Microbiol*. 2012; 20:209-210.

Gribaldo S, Poole AM, Daubin V, Forterre P, Brochier-Armanet C. The origin of eukaryotes and their relationship with the Archaea: are we at a phylogenomic impasse? *Nat Rev Microbiol*. 2010; 8:743-752.

Guy L, Ettema TJ. The archaeal 'TACK' superphylum and the origin of eukaryotes. *Trends Microbiol*. 2011; 19:580-587.

Jurgens G, Lindstrom K, Saano A. Novel group within the kingdom Crenarchaeota from boreal forest soil. *Appl Environ Microbiol*. 1997; 63:803–805.

Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol*. 2013; 30:772-780.

Kostka M, Uzlikova M, Cepicka I, Flegr J. SlowFaster, a user-friendly program for slow-fast analysis and its application on phylogeny of Blastocystis. *BMC Bioinformatics*. 2008; 9:341.

López-García P, Moreira D. Metabolic symbiosis at the origin of eukaryotes. *Trends Biochem Sci*. 1999; 24:88-93.

Moore CE, Archibald JM. Nucleomorph genomes. *Annu Rev Genet*. 2009; 43:251-264.

Narasimharao P, Podell S, Ugalde JA, Brochier-Armanet C, Emerson JB, Brocks JJ, Heidelberg KB, Banfield JF, Allen EE. De novo metagenomic assembly reveals abundant novel major lineage of Archaea in hypersaline microbial communities. *ISME J*. 2012; 6:81-93.

Nelson-Sathi S, Dagan T, Landan G, Janssen A, Steel M, McInerney JO, Deppenmeier U, Martin WF. Acquisition of 1,000 eubacterial genes physiologically transformed a methanogen at the origin of Haloarchaea. *Proc Natl Acad Sci U S A*. 2012; 109:20537-20542.

Nunoura T, Takaki Y, Kakuta J, Nishi S, Sugahara J, Kazama H, Chee GJ, Hattori M, Kanai A, Atomi H, Takai K, Takami H. Insights into the evolution of Archaea and eukaryotic protein modifier systems revealed by the genome of a novel archaeal group. *Nucleic Acids Res.* 2011; 39:3204-3223.

Philippe H. MUST, a computer package of Management Utilities for Sequences and Trees. *Nucleic Acids Res.* 1993; 21:5264-5272.

Philippe H, Brinkmann H, Lavrov DV, Littlewood DT, Manuel M, Wörheide G, Baurain D. Resolving difficult phylogenetic questions: why more sequences are not enough. *PLoS Biol.* 2011; 9:e1000602.

Philippe H, Laurent J. How good are deep phylogenetic trees? *Curr Opin Genet Dev.* 1998; 8:616-623.

Randau L, Münch R, Hohn MJ, Jahn D, Söll D. *Nanoarchaeum equitans* creates functional tRNAs from separate genes for their 5'- and 3'-halves. *Nature.* 2005; 433:537-541.

Rodríguez-Ezpeleta N, Brinkmann H, Roure B, Lartillot N, Lang BF, Philippe H. Detecting and overcoming systematic errors in genome-scale phylogenies. *Syst Biol.* 2007; 56:389-399.

Ronquist F, Teslenko M, van der Mark P, Ayres DL, Darling A, Höhna S, Larget B, Liu L, Suchard MA, Huelsenbeck JP. MrBayes 3.2: efficient Bayesian phylogenetic inference and model choice across a large model space. *Syst Biol.* 2012; 61:539-542.

Roure B, Baurain D, Philippe H. Impact of missing data on phylogenies inferred from empirical phylogenomic data sets. *Mol Biol Evol.* 2013; 30:197-214.

Schleper C, Holben W, Klenk HP. Recovery of crenarchaeotal ribosomal DNA sequences from freshwater-lake sediments. *Appl Environ Microbiol.* 1997; 63:321-323.

Spang A, Hatzenpichler R, Brochier-Armanet C, Rattei T, Tischler P, Spieck E, Streit W, Stahl DA, Wagner M, Schleper C. Distinct gene set in two different lineages of ammonia-oxidizing archaea supports the phylum Thaumarchaeota. *Trends Microbiol.* 2010; 18:331-340.

Stackebrandt E, Murray RGE, Trüper HG. *Proteobacteria* classis nov., a name for the phylogenetic taxon that includes the "purple bacteria and their relatives". *Int J Syst Bacteriol.* 1988; 38:321-325.

Stamatakis A. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models. *Bioinformatics.* 2006; 22:2688-2690.

Waters E, Hohn MJ, Ahel I, Graham DE, Adams MD, Barnstead M, Beeson KY, Bibbs L, Bolanos R, Keller M, Kretz K, Lin X, Mathur E, Ni J, Podar M, Richardson T, Sutton GG, Simon M, Soll D, Stetter KO, Short JM, Noordewier M. The genome of *Nanoarchaeum equitans*: insights into early archaeal evolution and derived parasitism. *Proc Natl Acad Sci U S A.* 2003; 100:12984-12988.

Table 1

Revised classification of Archaea into two major phyla and their respective classes and orders

Phylum	Class	Order	
Euryarchaeota	Archaeoglobi	Archaeoglobales	
	Methanobacteria	Methanobacteriales	
	Methanococci	Methanococcales	
	Methanomicrobia		Methanocellales
			Methanomicrobiales
			Methanosarcinales
			Halobacteriales
		Methanopyri	Methanopyrales
	Nanoarchaea	Nanoarchaeales	
	Thermococci	Thermococcales	
Thermoplasmata	Thermoplasmatales		
Proteoarchaeota	Aigarchaea	Aigarchaeales	
	Crenarchaea*		Acidilobales
			Desulfurococcales
			Fervidicoccales
			Sulfolobales
			Thermoproteales
	Korarchaea	Korarchaeales	
	Thaumarchaea		Cenarchaeales
			Nitrosopumilales
			Nitrososphaerales

* The current class Thermoprotei would be synonym of Crenarchaea.

Figure legends

FIG. 1.- Bayesian phylogenetic tree of Archaea rooted on bacterial sequences. The tree is based on the concatenation of 32 ribosomal proteins (2560 sites). Numbers at branches are posterior probabilities. The scale bar indicates the number of substitutions per position.

FIG. 2.- Bayesian phylogenetic tree of Archaea rooted on bacterial sequences. The tree is based on the concatenation of 38 diverse conserved proteins (6890 sites). Numbers at branches are posterior probabilities. The scale bar indicates the number of substitutions per position.

FIG. 3.- Bayesian phylogenetic tree of Archaea rooted on bacterial sequences. The tree is based on the concatenation of 32 ribosomal proteins and 38 diverse conserved proteins (9450 sites). Numbers at branches are posterior probabilities. The scale bar indicates the number of substitutions per position.

FIG. 4.- Unrooted Bayesian phylogenetic tree of Archaea including fast-evolving taxa. The tree is based on the concatenation of 32 ribosomal proteins and 38 diverse conserved proteins (9450 sites). Numbers at branches are posterior probabilities. The scale bar indicates the number of substitutions per position.

Supplementary material

Supplementary table S1. List of archaeal and bacterial complete genome sequences used in this work.

Supplementary table S2. List of all protein markers used in this work.

Supplementary table S3. List of species-pairs used for evolutionary estimations for gene-by-gene desaturation analysis.

Supplementary fig. S1. Maximum likelihood tree based on the complete set of markers (32 ribosomal proteins and 38 diverse conserved proteins, 9450 sites) and rooted using bacterial sequences as outgroup. Numbers at branches are bootstrap proportions. The scale bar indicates the number of substitutions per position.

Supplementary fig. S2. Unrooted maximum likelihood tree based on the complete set of markers (32 ribosomal proteins and 38 diverse conserved proteins, 9450 sites). Numbers at branches are bootstrap proportions. The scale bar indicates the number of substitutions per position.

Supplementary figs. S3-S12. Rooted maximum likelihood trees corresponding to the site-by-site desaturation analysis. Numbers at branches are bootstrap proportions. The scale bar indicates the number of substitutions per position.

Supplementary figs. S13-S22. Rooted maximum likelihood trees corresponding to the gene-by-gene desaturation analysis. Numbers at branches are bootstrap proportions. The scale bar indicates the number of substitutions per position.



Figure 1
101 species - 2560 positions 0.4

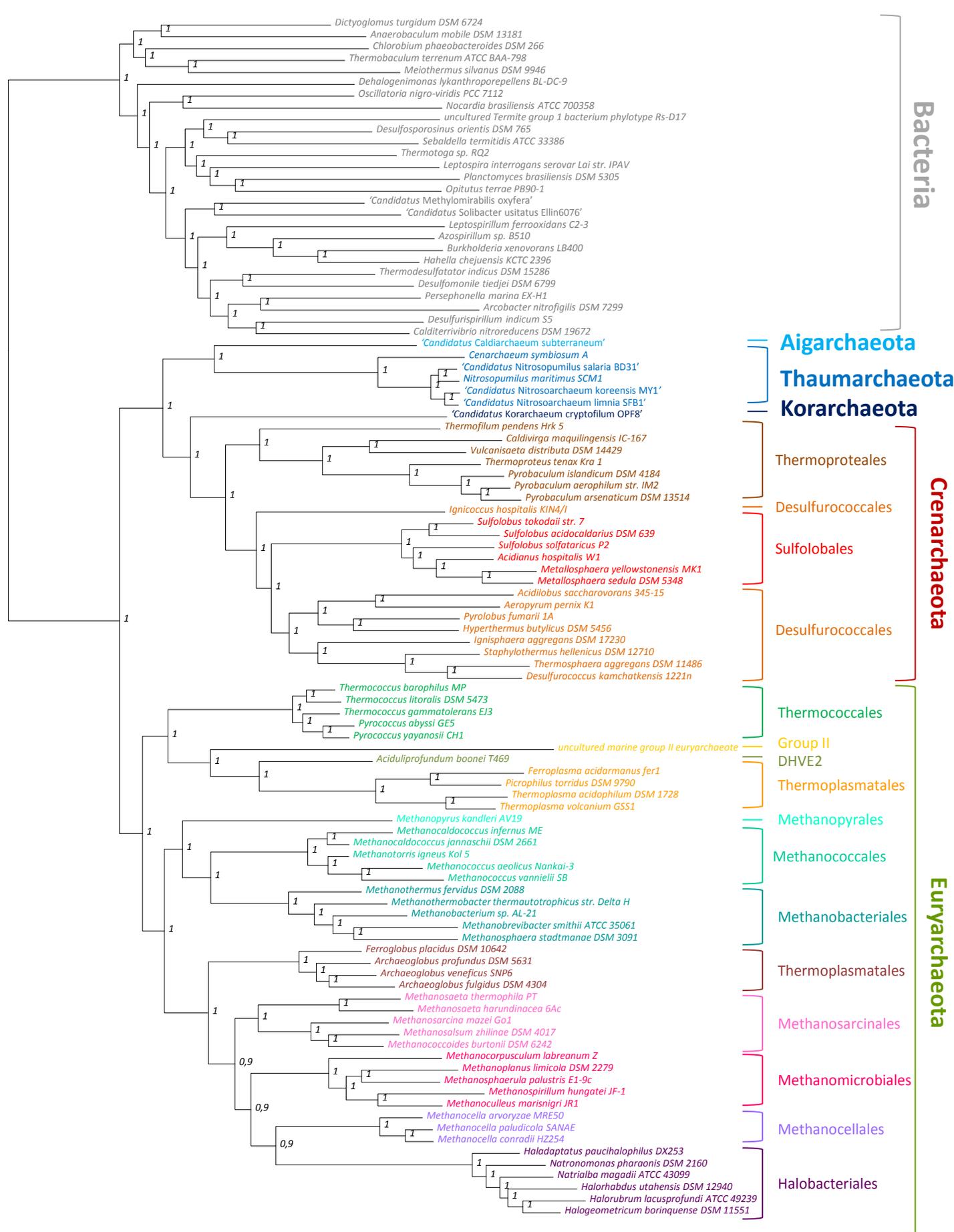


Figure 2
101 species - 9450 positions

0.2

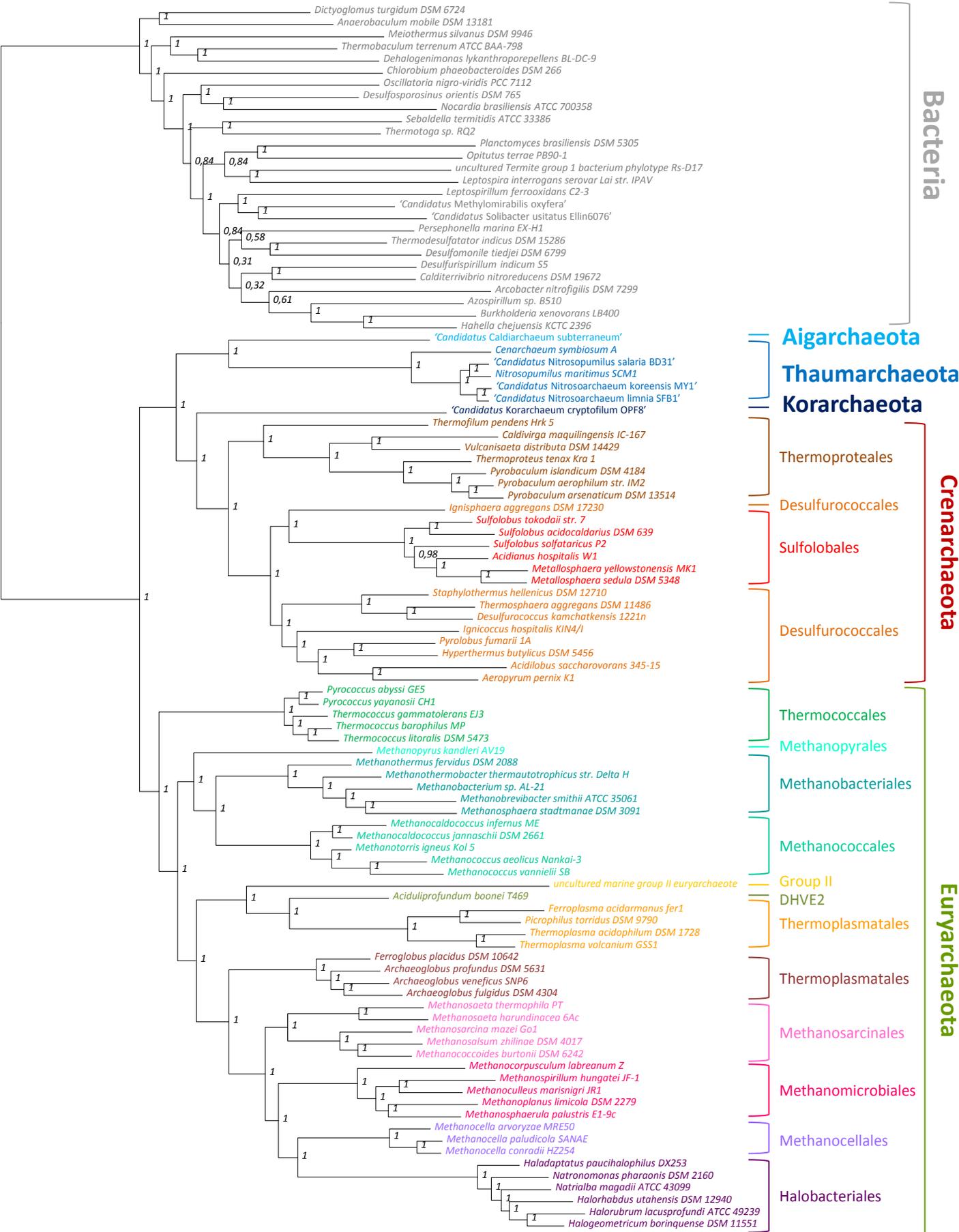


Figure 3
101 species - 9540 positions

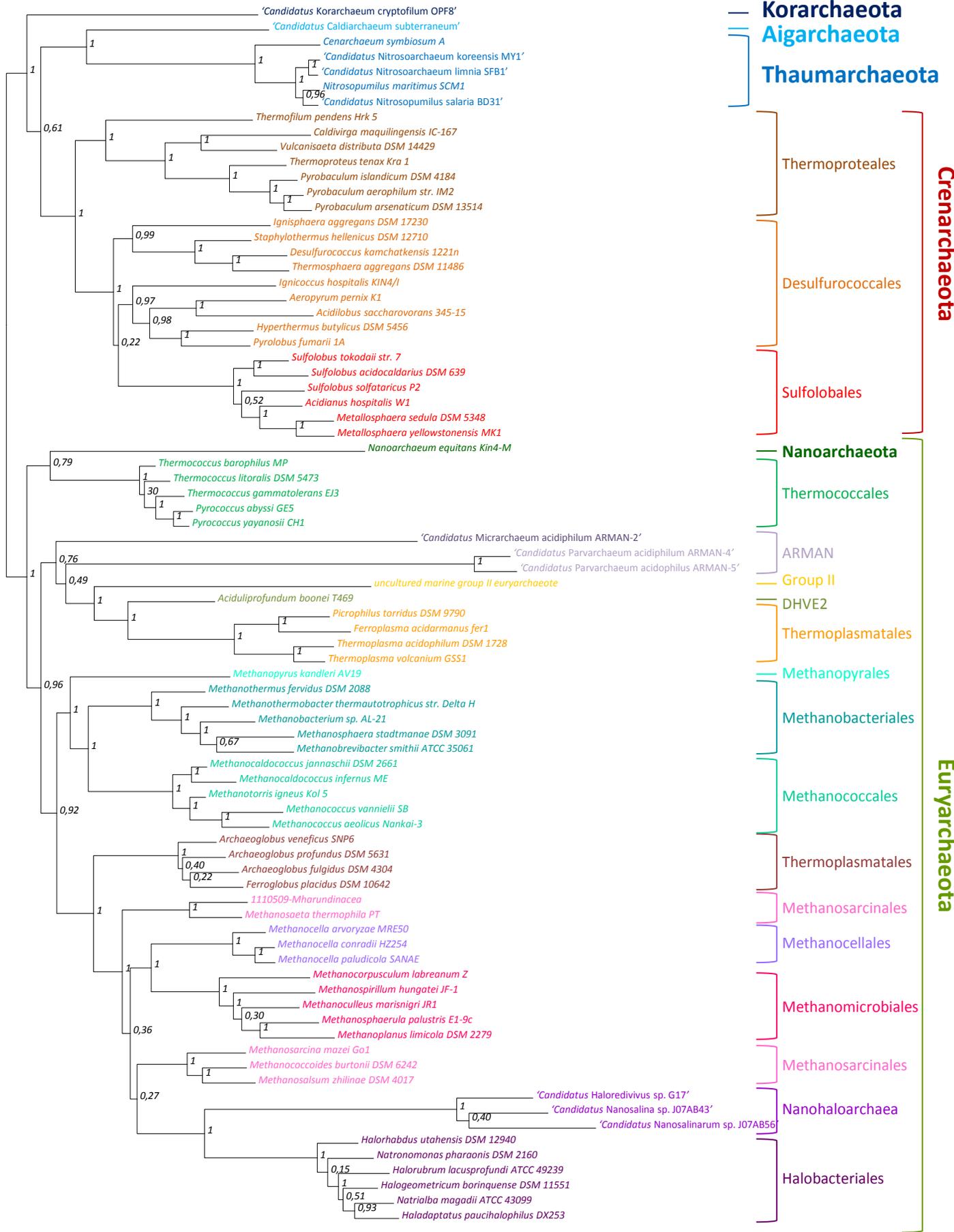


Figure 4
81 species - 9540 positions

0.2

C. Synthèse et éléments de discussion

Dans cet article, nous avons donc cherché à placer la racine de l'arbre des archées à partir de deux sets de protéines, 32 protéines ribosomiques et 38 nouveaux marqueurs impliqués dans différents processus cellulaires, avec leur homologues bactériens.

Nous avons tout d'abord écarté de nos premières analyses les lignées problématiques des genres *Nanoarchaeum*, '*Ca. Micrarchaeum*', '*Ca. Parvarchaeum*' et les Nanohaloarchaeota, afin de diminuer le risque de LBA. La première phylogénie inférée à partir des 32 protéines ribosomiques en maximum de vraisemblance et en bayésien place la racine au sein des Euryarchaeota. Différents groupes d'euryarchées émergent de façon paraphylétique et avec un très faible soutien, alors que les Thaumarchaeota, Aigarchaeota, Crenarchaeota et Korarchaeota (le superphylum « TACK » proposé par Guy et Ettema en 2011 (Guy and Ettema 2011)) groupent de façon monophylétique bien soutenue à l'intérieur des euryarchées. Les analyses du second jeu de données, contenant 38 nouveaux marqueurs, placent la racine entre les Euryarchaeota, monophylétiques, et le groupe « TACK » avec un support au nœud maximal. La concaténation de ces deux jeux de données permet aussi d'obtenir une phylogénie racinée entre les Euryarchaeota et le groupe « TACK » très bien soutenue. A partir de ce résultat, l'analyse en bayésien de la position des archées de petite taille dans cette phylogénie mais sans groupe extérieur (pour minimiser les problèmes de LBA) a montré qu'elles se replaçaient au sein des euryarchées et n'étaient pas groupées. La racine de l'arbre des archées serait donc entre les Euryarchaeota et le superphylum « TACK ».

Une diversité phylogénétique et physiologique d'une ampleur similaire à celle rencontrée au sein des Euryarchaeota (qui leur a valu leur nom, *eurus* voulant dire ample, spacieux) se retrouve aussi dans le groupe « TACK ». Cette observation, ainsi que la monophylie très robuste de ce groupe et la position de la racine le séparant des Euryarchaeota, nous ont conduits à proposer un changement majeur dans la taxonomie des archées, à savoir, la création d'un nouveau phylum 'Proteoarchaeota' afin de donner un nom formel au groupe « TACK ». Ceci doit être accompagné du déclassement des phyla composant le groupe « TACK » en ordres : Thaumarchaeales, Aigarchaeales, et Korarchaeales ; les Crenarchaeota deviendraient les Crenarchaea. Nous proposons aussi de supprimer le phylum Nanoarchaeota et de le considérer comme un ordre au sein des Euryarchaeota sous le nom de Nanoarchaeales.

La multiplication des phyla proposés ces dernières années est assez symptomatique de l'absence de règles claires en ce qui concerne la taxonomie de haut niveau chez les archées et les bactéries, créant un flou dans l'appellation des différents groupes, et par conséquent dans la compréhension et la perception de leur phylogénie. Certains de ces phyla ont été proposés sur

l'analyse de très peu de données en termes de génomes séquencés appartenant au groupe et du nombre de marqueurs utilisés. La publication très récente de Rinke et collaborateurs 2013 (Rinke et al. 2013), décrivant le séquençage de 201 nouveaux génomes microbiens et proposant plusieurs nouveaux phyla d'archées, est un exemple de plus de cette mouvance. Or, notre travail montre qu'une analyse fine permet d'obtenir une phylogénie des archées racinée et de bonne qualité. Il est important de remarquer ici que ce résultat a été obtenu grâce à l'analyse d'autres marqueurs que les protéines ribosomiques, montrant que malgré la puissance de ces marqueurs pour la résolution de la phylogénie des archées, il peut être nécessaire d'utiliser des protéines différentes pour répondre à des questions différentes. Une analyse fine avec un jeu de données adapté à la question posée, en l'occurrence, les relations phylogénétiques anciennes entre les lignées d'archées, semble donc nécessaire avant la proposition d'un nouveau groupe de haut niveau taxonomique tel qu'un phylum. La nouvelle classification que nous proposons rétablit un équilibre entre les ordres d'Euryarchaeota et les groupes composant les 'Proteoarchaeota', en termes de diversité et de distance phylogénétique.

Le positionnement de la racine de l'arbre des archées est aussi très important en ce qui concerne la nature du dernier ancêtre commun de ce domaine. D'après notre analyse, il semblerait avoir été hyperthermophile, les premiers ordres divergeant au sein des Euryarchaeota et des 'Proteoarchaeota' étant hyperthermophiles ou très thermophiles (les Thermococcales et Nanoarchaeales chez les euryarchées et les Korarchaeales et Aigarchaeales chez les 'proteoarchées'). Cet ancêtre aurait été aussi probablement très complexe, contenant au moins les caractères partagés entre les deux grands phyla, et non pas extrêmement réduit à l'image des nanoarchées comme *N. equitans* ou les ARMAN.

Enfin, en ce qui concerne la relation des eucaryotes avec les archées, il serait intéressant de construire un nouveau set de marqueurs sélectionnés pour répondre à cette question précise. Les homologues bactériens n'étant plus nécessaires, ces marqueurs pourront être choisis indépendamment de la présence d'homologues bactériens, de façon à s'affranchir de la nécessité d'avoir des protéines universelles, ainsi que du risque de biais liés à la différence de vitesse d'évolution entre les domaines, ou de la saturation. Cette méthode a déjà été proposée par Gribaldo et collaborateurs en 2010 (Simonetta Gribaldo et al. 2010).

Chapitre 3 : Répartition taxonomique et histoire évolutive des protéines d'archées : exemple du système chaperonne DnaK et de la protéine DnaJ-Fer.

A. Introduction

Notre recherche de nouveaux marqueurs pour l'inférence de la phylogénie des archées est passée par l'analyse de la phylogénie de chacune des protéines codées dans les trois génomes de *Cenarchaeum symbiosum*, *Nitrosopumilus maritimus* et '*Ca. Caldiarchaeum subterraneum*'. Cette analyse non automatique a permis de noter un certain nombre d'observations sur ces phylogénies, telles que l'absence ou la présence d'homologues bactériens et eucaryotes, les phyla d'archées représentés, des cas de transferts horizontaux potentiels, ou des cas de paralogie. Ces données n'ont pas été exploitées dans ma thèse, le but premier de cette étude étant la recherche de marqueurs pour l'étude de la phylogénie des archées, mais elles seraient très intéressantes à exploiter dans le futur. Deux exemples en sont la preuve. Tout d'abord, le travail présenté dans le Chapitre 2 a été entamé car il a été possible de vérifier très rapidement si des homologues bactériens étaient présents dans les jeux de données de départ des 200 marqueurs sélectionnés pour l'inférence de la phylogénie des archées. L'autre exemple de l'intérêt de ce travail est de mettre en évidence des histoires évolutives particulières comme celle de la protéine DnaJ-Fer présentée dans l'article ci-dessous.

Cette protéine a une répartition taxonomique très particulière, elle n'est présente que chez certaines Thaumarchaeota et les Viridiplantae (algues et plantes vertes), ce qui a attiré notre attention. L'étude de l'histoire évolutive de cette protéine nous a amenés à nous intéresser aussi à l'histoire évolutive du système chaperonne DnaK. En effet, les protéines contenant un domaine peptidique dnaJ, sont connues pour interagir avec DnaK dans les réponses aux chocs thermiques et d'autres types de stress. La protéine DnaK et ses co-chaperonnes GrpE et DnaJ sont présentes de façon ubiquitaire chez les bactéries et les eucaryotes, mais sont absentes chez certaines archées, et totalement absentes chez les archées hyperthermophiles. Il a été montré que la protéine DnaJ-Fer, composée d'un domaine dnaJ et d'un domaine ferredoxine, intervient en association avec DnaK dans le chloroplaste de l'algue verte *Chlamydomonas reinhardtii* (Dorn et al. 2010). Nous avons montré que chez les archées, DnaK, GrpE et DnaJ (une protéine différente de la protéine DnaJ-Fer, contenant aussi un domaine DnaJ et connue pour interagir avec DnaK) avaient été acquises par transfert horizontal de gène depuis des bactéries. Au moins deux évènements indépendants de transferts ont été observés, un vers les Halobacteriales et les Nanoarchaea et un autre vers d'autres euryarchées. Elles ont ensuite été transférées plusieurs fois au sein des archées. C'est le cas

chez les Thaumarchaeota, qui auraient acquis ce système depuis d'autres archées (euryarchées), même si le donneur exact n'a pas pu être identifié. Parallèlement, nous avons aussi mis en évidence que la protéine DnaJ-Fer était issue, chez les Thaumarchaeota, d'un transfert depuis les Viridiplantae, celles-ci ayant acquis au moins le domaine ferredoxine depuis des cyanobactéries lors de l'endosymbiose chloroplastique. Le domaine dnaJ ne semble pas avoir la même origine chez les Viridiplantae et chez les Thaumarchaeota mais l'origine commune du domaine ferredoxine nous incite à penser que la protéine a été transférée déjà fusionnée et qu'un remplacement homologue entre deux domaines dnaJ s'est produit chez les thaumarchées. La répartition taxonomique du système DnaK-GrpE-DnaJ et DnaJ-Fer et leur fonction cellulaire nous ont amenés à proposer que ces quatre protéines étaient peut-être impliquées dans l'adaptation à la mésophilie chez les archées et particulièrement chez les thaumarchées.

Cette étude a fait l'objet de la publication « *Horizontal gene transfer of a chloroplast DnaJ-Fer protein to Thaumarchaeota and the evolutionary history of the DnaK chaperone system in Archaea* » (Petitjean, Moreira, López-García, Brochier-Armanet) dans le journal *BMC Evolutionary Biology* (Petitjean et al. 2012).

Chapitre 3 : Répartition taxonomique et histoire évolutive des protéines d'archées : exemple du système chaperonne DnaK et de la protéine DnaJ-Fer.

B. Manuscrit de l'article 3 : « Horizontal gene transfer of a chloroplast DnaJ-Fer protein to Thaumarchaeota and the evolutionary history of the DnaK chaperone system in Archaea »

RESEARCH ARTICLE

Open Access

Horizontal gene transfer of a chloroplast DnaJ-Fer protein to Thaumarchaeota and the evolutionary history of the DnaK chaperone system in Archaea

Céline Petitjean^{1,2}, David Moreira², Purificación López-García² and Céline Brochier-Armanet^{3*}

Abstract

Background: In 2004, we discovered an atypical protein in metagenomic data from marine thaumarchaeotal species. This protein, referred as DnaJ-Fer, is composed of a J domain fused to a Ferredoxin (Fer) domain. Surprisingly, the same protein was also found in Viridiplantae (green algae and land plants). Because J domain-containing proteins are known to interact with the major chaperone DnaK/Hsp70, this suggested that a DnaK protein was present in Thaumarchaeota. DnaK/Hsp70, its co-chaperone DnaJ and the nucleotide exchange factor GrpE are involved, among others, in heat shocks and heavy metal cellular stress responses.

Results: Using phylogenomic approaches we have investigated the evolutionary history of the DnaJ-Fer protein and of interacting proteins DnaK, DnaJ and GrpE in Thaumarchaeota. These proteins have very complex histories, involving several inter-domain horizontal gene transfers (HGTs) to explain the contemporary distribution of these proteins in archaea. These transfers include one from Cyanobacteria to Viridiplantae and one from Viridiplantae to Thaumarchaeota for the DnaJ-Fer protein, as well as independent HGTs from Bacteria to mesophilic archaea for the DnaK/DnaJ/GrpE system, followed by HGTs among mesophilic and thermophilic archaea.

Conclusions: We highlight the chimerical origin of the set of proteins DnaK, DnaJ, GrpE and DnaJ-Fer in Thaumarchaeota and suggest that the HGT of these proteins has played an important role in the adaptation of several archaeal groups to mesophilic and thermophilic environments from hyperthermophilic ancestors. Finally, the evolutionary history of DnaJ-Fer provides information useful for the relative dating of the diversification of Archaeplastida and Thaumarchaeota.

Keywords: DnaJ/Hsp40, DnaK/Hsp70, Hyperthermophily, Archaeplastida, Phylogeny, Archaea, Thaumarchaeota, Horizontal gene transfer, Mesophily

Background

The 70 kD heat shock proteins (called DnaK in bacteria and Hsp70 in eukaryotes) form a large family of molecular chaperones upregulated in cells suffering various stresses, including heat shocks and heavy metal exposure [1,2]. In addition, these proteins play a major role during protein synthesis by binding to the nascent peptides exiting the ribosome in order to prevent their aggregation and facilitating their folding in the optimal functional conformation [3]. During the interaction with the

partially synthesized peptides, DnaK/Hsp70 increases its ATPase activity [3]. This chaperone has two main partners: the J-proteins [4,5] and the nucleotide exchange factor, called GrpE in bacteria (or Mge1 [6] in mitochondria and Cge1 [7] in chloroplasts) and Bag-1, a eukaryotic functional analogue of GrpE [8]. The nucleotide exchange factor promotes the exchange of ADP to fresh ATP in the nucleotide-binding region of DnaK/Hsp70, whereas the J-proteins stimulate the ATPase activity in order to stabilize the interaction of DnaK with unfolded proteins [5,9,10]. The J-proteins form a large family of proteins, which are structurally and functionally diverse but all have the capacity to interact with DnaK/Hsp70 through their J-domain [4,11]. Among them, DnaJ/Hsp40 proteins form the largest subfamily [12]. They control the flux of unfolded

* Correspondence: celine.brochier-armanet@univ-lyon1.fr

³CNRS, UMR5558, Laboratoire de Biométrie et Biologie Evolutive, Université de Lyon, Université Lyon 1, 43 boulevard du 11 novembre 1918, 69622, Villeurbanne, France

Full list of author information is available at the end of the article

polypeptides into and out of the substrate-binding domain of DnaK/Hsp70 [9,11].

DnaK proteins are widespread, being encoded by a single gene in most bacterial genomes, whereas most eukaryotic genomes harbor several Hsp70 genes that may have diverse evolutionary origins [1,13,14]. For example, in the green alga *Chlamydomonas reinhardtii*, five Hsp70 copies are present, all them encoded in the nuclear genome despite being targeted in diverse cellular compartments: three of them most likely originated by duplications from an ancestral eukaryotic gene (one expressed in the cytoplasm and two in the endoplasmic reticulum); one has a mitochondrial origin and is exported into the mitochondria, whereas the latter originated from the chloroplast endosymbiosis and is targeted into the chloroplast [15]. In contrast with DnaK, the J-proteins are encoded in multiple copies in bacterial genomes [9]. This is also the case in eukaryotes, where they work in the different cell compartments in association with the Hsp70 proteins cited above [9,11]. Finally, the nucleotide exchange factor GrpE is present in one copy in most of bacterial genomes, whereas the eukaryotic Mge1, Cge1 and Bag-1 are encoded in the nucleus but addressed to the mitochondria, chloroplasts, and to the nucleus and the cytoplasm, respectively [7,8].

The presence of DnaK, DnaJ and GrpE has been reported in several archaeal genomes [16], more precisely in several euryarchaeota but never in crenarchaeotal species. The best studied case concerns DnaK. A phylogenetic analysis by Gribaldo and coworkers suggested that this protein was acquired by several archaea by horizontal gene transfer (HGT) from different bacterial donors [17]. These authors observed three different groups of archaeal DnaK sequences branching specifically with certain bacterial homologues. More precisely, *Methanosarcina mazei* (Methanosarcinales) was related to the *Clostridium* group of Firmicutes (low G+C Gram positive bacteria), *Halobacterium cutirubrum* and *Halobacterium marismortui* (Halobacteriales) to the Actinobacteria (high G+C Gram positive bacteria), whereas *Methanobacterium thermautotrophicum* (Methanobacteriales) and *Thermoplasma acidophilum* (Thermoplasmatales) branched with *Thermotoga maritima* (Thermotogales) [17]. More recently, Macario et al. (2006) studied in various bacteria and archaea the taxonomic distribution and the phylogeny not only of DnaK but also of GrpE and DnaJ. They showed that the genes coding for these three proteins were clustered in most of the genomes examined [16]. They also confirmed the results of Gribaldo et al. (1999), i.e. the likely existence of three HGT events from bacteria to archaea. However, they proposed a more complex scenario where the DnaK/DnaJ/GrpE cluster was first acquired

from a bacterial donor by the ancestor of the Euryarchaeota, then lost in Methanococcales and in the common ancestor of Archaeoglobales, Halobacteriales and Methanosarcinales, and finally reacquired independently by Halobacteriales and Methanosarcinales from Actinobacteria and from Firmicutes, respectively [16]. Worth noting, in these two studies, none of the three proteins was detected in hyperthermophilic archaea.

In addition to these relatively well-characterized chaperones and co-chaperones, the study of a genomic fragment of an uncultured deep marine archaeon from an environmental DNA fosmid library revealed a very unusual J-protein, referred as DnaJ-Fer, composed of a J-domain fused with a Ferredoxin (Fer) domain [18]. The phylogenetic analysis of a 16S rRNA gene also found in this genomic fragment showed that it belonged to a member of the Thaumarchaeota, more precisely in the I.1a subgroup. These archaea, formerly classified as Group I, a sublineage of Crenarchaeota [19,20], have been recently proposed to represent a third phylum of Archaea together with the Euryarchaeota and Crenarchaeota [21]. Thaumarchaeota are widespread in many environments, including marine and freshwater, soil and sediment [22,23]. Surprisingly, the presence of DnaJ-Fer proteins has also been reported in Viridiplantae (including green algae and plants), with three homologues (CDJ3, 4 and 5) in *C. reinhardtii* [24]. These proteins are localized in the chloroplast of this green alga where they interact with the chloroplast Hsp70B and Cge1 proteins. However, the precise function of these DnaJ-Fer proteins in *C. reinhardtii* remains to be elucidated. According to the location and the nature of its partners, it would be tempting to hypothesize a cyanobacterial origin of the DnaJ-Fer protein. However, no homologue has been detected in Cyanobacteria [24].

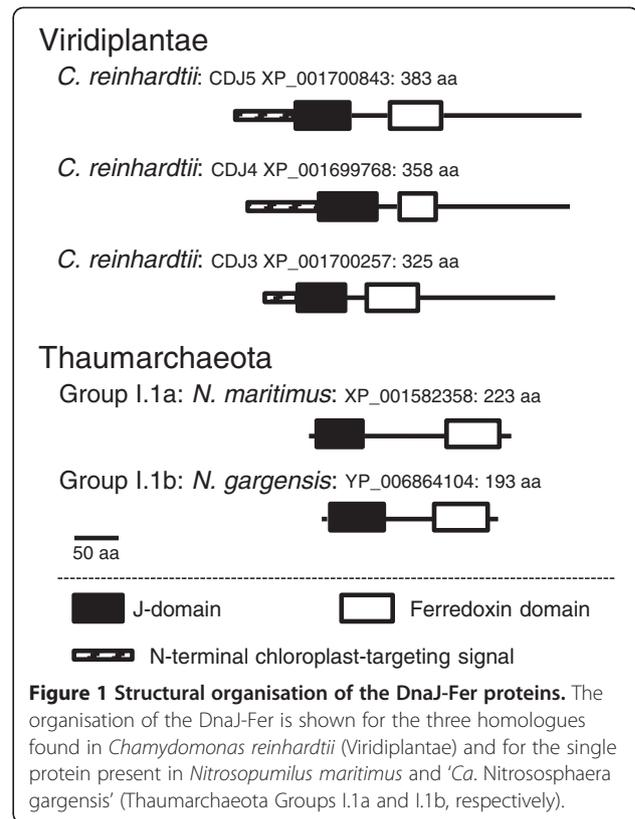
Two hypotheses can explain the unexpected taxonomic distribution of the DnaJ-Fer protein in Thaumarchaeota and Viridiplantae: either two independent and convergent fusions of the two protein domains occurred in these two distantly related lineages, or a single fusion occurred in one of them followed by a HGT to the other lineage [24]. In this work, we have taken advantage of the recent burst of available archaeal complete genome sequences [25], including representatives of new major lineages such as the Thaumarchaeota, ARMAN or Nanohaloarchaeales, to decipher the evolutionary history of DnaK and its co-chaperones in Archaea, with especial attention on the intriguing DnaJ-Fer protein. Our results support a complex scenario in which HGT appears to have played an important role. In addition to other cases of HGT, Thaumarchaeota appear to have most likely acquired their DnaK, co-chaperones and DnaJ-Fer proteins by independent HGTs from multiple donors, including other archaea and plants.

Results

DnaJ-Fer proteins are widespread in viridiplantae and thaumarchaeota

We carried out an intensive survey of public sequence databases to find that DnaJ-Fer homologues are present in all Viridiplantae (green algae and land plants) for which complete genome sequences were available. In contrast, we did not detect them in Rhodophyta and Glaucophyta, the two other lineages composing the Plantae or Archaeplastida eukaryotic supergroup [26]. However, due to the scarcity of sequence data from these two lineages, we can not exclude the future discovery of DnaJ-Fer in some species belonging to them. In addition to green algae and land plants, DnaJ-Fer homologues were detected in the four available complete genomes of Thaumarchaeota (Additional file 1): *Cenarchaeum symbiosum* (a sponge symbiont) [27], the planktonic *Nitrosopumilus maritimus* (the first isolated thaumarchaeote) [28] and its two close relatives ‘*Candidatus* (*Ca.*) Nitrosoarchaeum limnia SFB1’ [29] and in ‘*Ca.* Nitrosoarchaeum koreensis MY1’ [30] which live in low salinity sediments and in the soil rhizosphere, respectively, as well as in several environmental fosmid sequences, all likely members of the mesophilic group I.1a. The protein was also present in *Nitrososphaera viennensis* (Schleper and Spang, personal communication) and ‘*Ca.* Nitrososphaera gargensis’ [31], two moderate thermophilic representatives of the group I.1b. In contrast, it was absent in the thermophilic species ‘*Ca.* Nitrosocaldus yellowstonii’, a representative of the more distant Hot Water Crenarchaeotic Group (HWCG) III (de la Torre, personal communication), and in ‘*Ca.* Caldiarchaeum subterraneum’ [32], a representative of the ‘Aigarchaeota’ (formerly group HWCG I) which seems to be either the sister group of Thaumarchaeota or a deeply branching thaumarchaeotal lineage [22].

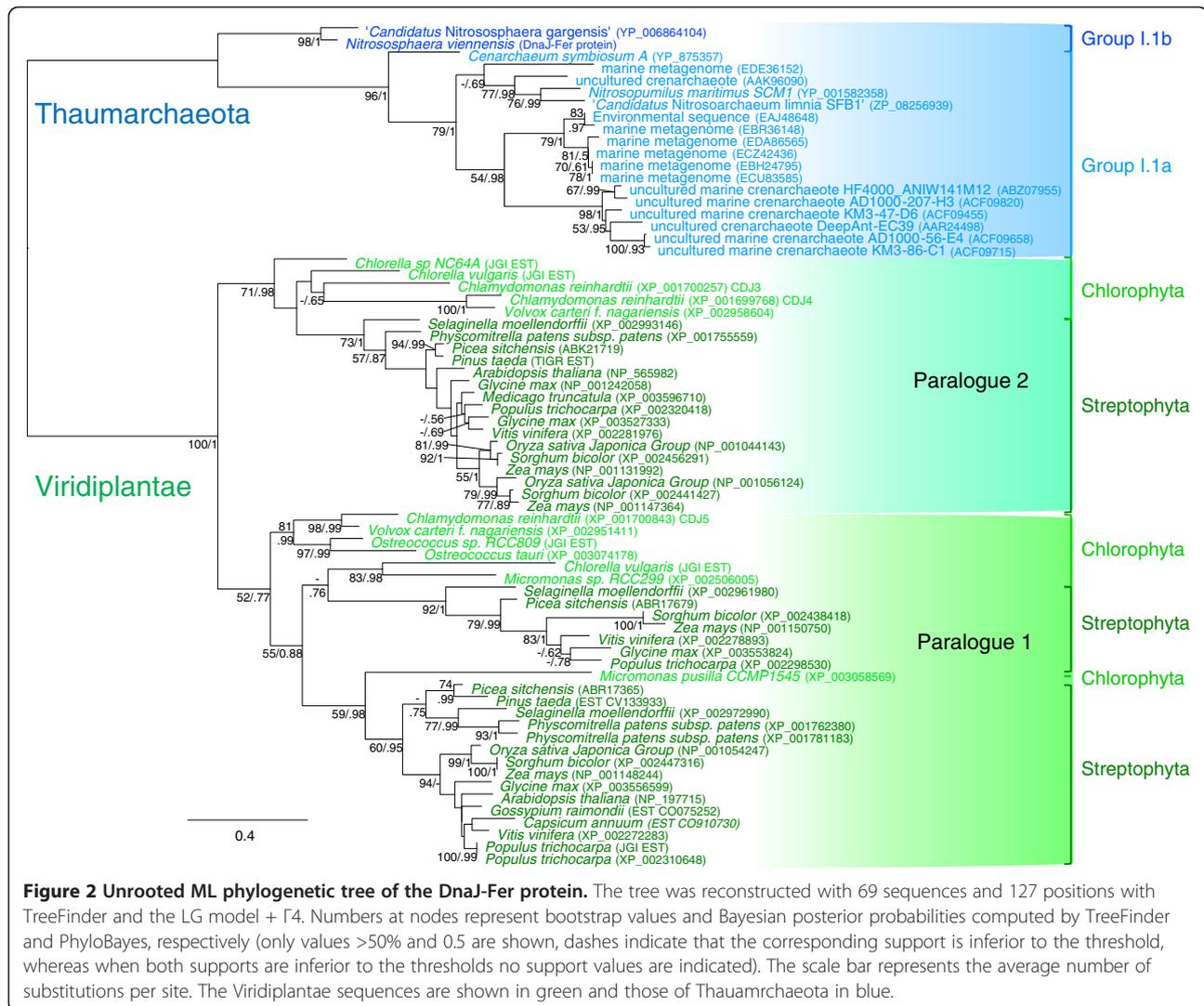
The J and Fer domains are two small domains (less than 100 amino acids) well conserved in the plant and thaumarchaeotal DnaJ-Fer sequences. The Fer domain was characterized by an amino acid motif CXXCXXC observed in all those sequences except in ‘*Ca.* *N. gargensis*’ and *N. viennensis*, where the motif was CXXFXXC. Contrasting with the conservation of these two domains, we observed different sequence organizations of the DnaJ-Fer proteins in the Viridiplantae and in the Thaumarchaeota (Figure 1). In Viridiplantae, an N-terminal chloroplast signal region preceded the J and the Fer domains, and the protein ended with a long C-terminal region (up to 150 amino acids) of unknown function. In Thaumarchaeota, these N- and C-terminal regions were absent, but an inter-domain region (ranging between 54 and 92 amino acids) was present between the J and the Fer domains. This region was well conserved in *N. maritimus*, *C. symbiosum*, ‘*Ca.* Nitrosoarchaeum limnia SFB1’,



‘*Ca.* Nitrosoarchaeum koreensis MY1’ and the fosmids found in the environmental database (all belonging to the group I.1a), but was divergent and shorter (54 amino acids) in the sequences of ‘*N. gargensis*’ and *N. viennensis*, the two representatives of group I.1b. The presence of this variable central region suggested that its role is probably structural and not functional in Thaumarchaeota. By contrast, much shorter or no central regions were present between the two domains in the plant sequences.

The taxonomic distribution of the DnaJ-Fer protein results from an ancient HGT

Maximum likelihood (ML) analyses and Bayesian inference (BI) of the DnaJ-Fer alignment revealed three monophyletic groups (Figure 2). Two corresponded to Viridiplantae (ML bootstrap values (BV) = 71% and 52%, and BI posterior probabilities (PP) = 0.98 and 0.78, respectively) whereas the third gathered the thaumarchaeotal sequences (BV = 100% and PP = 1.00). Interestingly, the relationships among sequences within each of these groups were in agreement with the accepted species phylogeny and relatively well supported despite the small number positions (127 amino acids) kept for the phylogenetic analysis. More precisely, the dichotomy between group I.1a and group I.1b Thaumarchaeota was well supported (BV = 96% and PP = 1.00). The relationships among the green algae and land plant



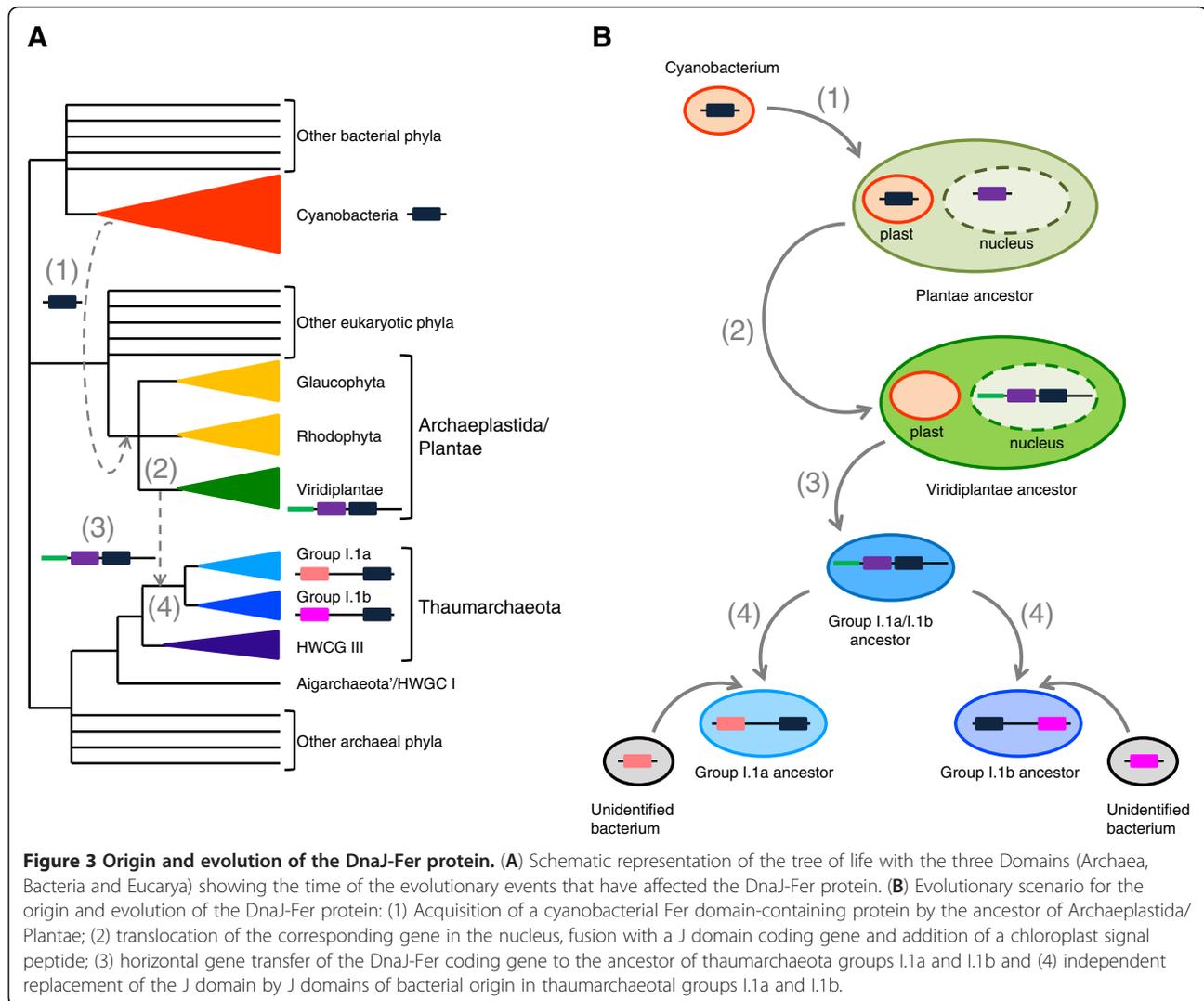
sequences were more complex since there were several copies of this protein in these species, most likely resulting from duplication events. A first duplication occurred almost certainly in the ancestor of Viridiplantae, leading to the two paralogues present in green algae and land plants. This event was followed by additional duplication events at the origin of the multiple copies of paralogues 1 and 2 observed in the viridiplantae lineages (Figure 2).

Phylogeny results indicated that the ancestor of thaumarchaeotal groups I.1a and I.1b already harboured the DnaJ-Fer gene and that the ancestor of Viridiplantae had two copies. If the unusual taxonomic distribution of DnaJ-Fer proteins was indicative of an HGT between Thaumarchaeota and Viridiplantae, the inferred phylogenies suggested that this HGT took place before the diversification of these two major lineages and was therefore relatively ancient (event 3 on Figure 3A). However, due to the lack of any suitable outgroup (no other lineage contained the DnaJ-Fer protein) it was not possible

to determine the precise evolutionary origin of the DnaJ-Fer gene and the direction of the HGT between Thaumarchaeota and Viridiplantae. To tackle this issue we carried out phylogenetic analyses of the J and Fer domains separately. Indeed, although the association between these two domains is specific of Thaumarchaeota and Viridiplantae, each domain is widely distributed in present day organisms, opening the possibility to reconstruct rooted phylogenies for each of them.

The J and Fer domains have two different evolutionary origins

As expected because of the small number of conserved sequence positions, the ML phylogeny of the Fer domain was largely unresolved (data not shown). Nevertheless, the Fer domain of the DnaJ-Fer proteins of Viridiplantae and Thaumarchaeota branched within a single cluster, which also contained various bacterial and archaeal sequences. To improve the resolution of the phylogenetic



relationships between these sequences, we carried out an analysis of the sequences composing this cluster and close relatives using several more distantly related sequences as outgroup. The resulting ML tree supported the grouping of thaumarchaeotal and viridiplantae sequences (BV = 77% and PP = 0.99, Figure 4A), indicating that the Fer domain of the DnaJ-Fer proteins had a single origin and, most likely, that a HGT event occurred between these two distant lineages. Interestingly, Fer domains from cyanobacterial and stramenopile species branched in the same cluster (Figure 4A). Stramenopiles are eukaryotes that acquired a chloroplast secondarily from Rhodophyta [33]. Therefore, the grouping of viridiplantae, stramenopile and cyanobacterial sequences strongly suggested a cyanobacterial origin of the Fer domain in these two eukaryotic photosynthetic lineages, even if the sequences of the photosynthetic eukaryotes did not appear nested within the cyanobacterial sequences. In fact, this was likely due to a poor resolution of the phylogenetic tree, which is frequent in similar studies

of proteins of cyanobacterial origin, where most often only a sister-grouping of cyanobacteria and plant sequences is observed in phylogenetic trees [34]. The hypothesis of an HGT from plants to cyanobacteria can be discarded because the protein is present in *Gloeobacter*, which is a deeply branching cyanobacterial lineage that has diverged before the chloroplastic endosymbiosis and, consequently, before the origin of plants [35]. The HGT of the Fer domain from cyanobacteria to plants is also strongly supported by the functional data showing that the DnaJ-Fer protein is targeted to the chloroplast in the green alga *Chlamydomonas* [24]. It is important to notice that, in contrast with the two-domain DnaJ-Fer proteins of Viridiplantae and Thaumarchaeota, the stramenopile and cyanobacterial proteins were composed uniquely of the Fer domain. Thus, the association between the J and the Fer domains probably occurred in the Viridiplantae lineage after the divergence of the present-day three main Archaeplastida phyla (i.e., Viridiplantae, Rhodophyta and Glaucophyta) but prior to the

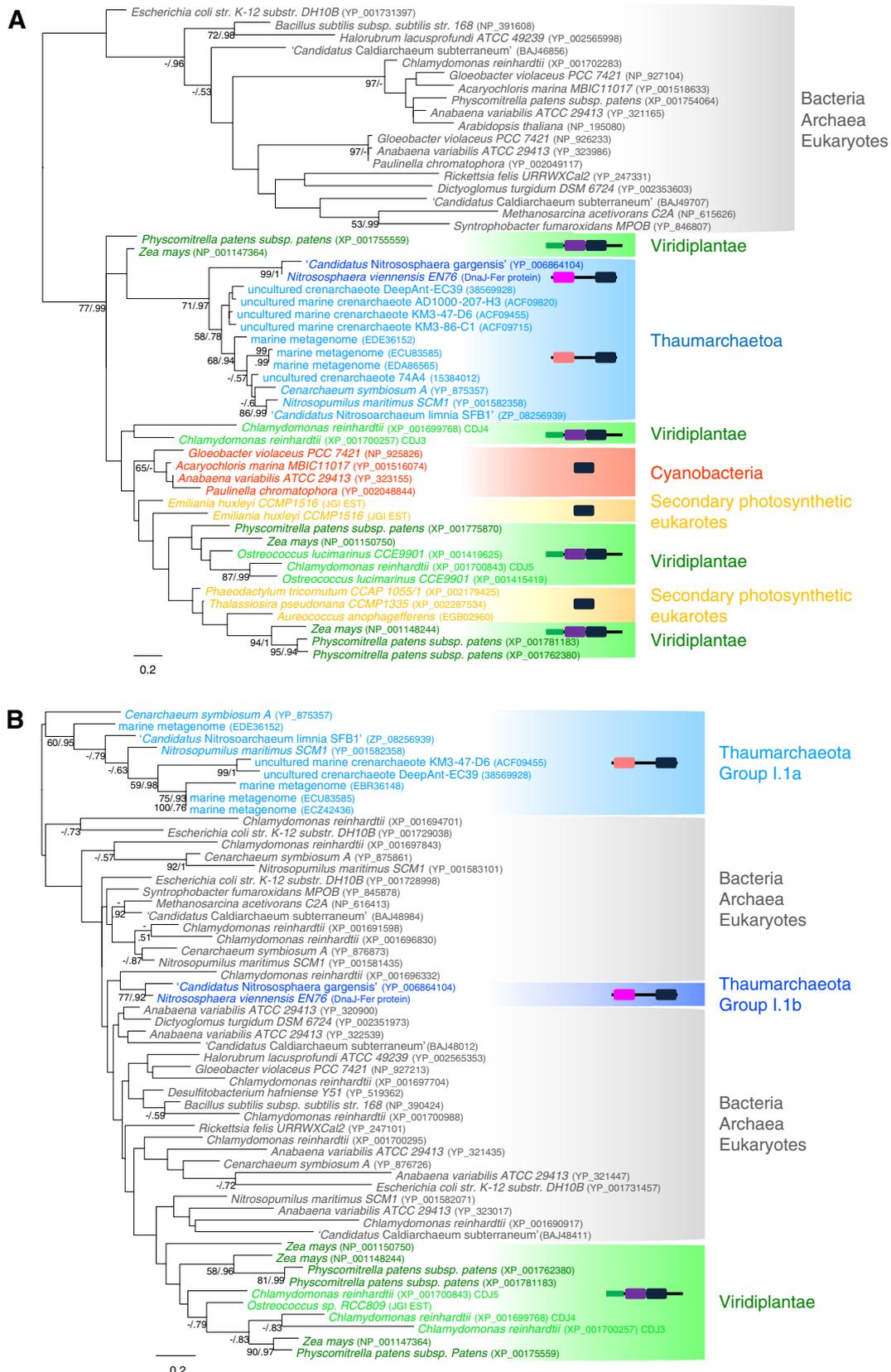


Figure 4 (See legend on next page.)

(See figure on previous page.)

Figure 4 Unrooted ML trees of the Fer and J domains. The ML tree of the Fer domain (A) was inferred with 52 sequences and 41 positions, whereas 55 sequences and 40 positions were kept to reconstruct the J domain tree (B). The two trees were inferred with TreeFinder (LG model). Numbers at nodes represent bootstrap values and Bayesian posterior probabilities computed with TreeFinder and PhyloBayes, respectively (only values >50% and 0.5 are shown, dashes indicate that the corresponding support is inferior to the threshold, whereas when both supports are inferior to the thresholds no support values are indicated). The scale bars represent the average number of substitutions per site. For clarity, the sequences relevant for the understanding of the history of DnaJ-Fer proteins have been coloured according to their taxonomy.

internal diversification of Viridiplantae (Figure 3A). This phylogeny also supported that the ancestor of the thaumarchaeotal groups I.1a and I.1b acquired secondarily the DnaJ-Fer protein from an ancestor of present-day Viridiplantae (Figure 3A). Another possibility would be that Viridiplantae and Cyanobacteria acquired their Fer domain from Thaumarchaeota. This would imply two HGT events, one from Thaumarchaeota to Cyanobacteria and a second one from Cyanobacteria to photosynthetic eukaryotes through the chloroplast endosymbiosis. In addition, that hypothesis would also imply the dissociation of the Fer and J domains in Cyanobacteria and their reassociation in the Viridiplantae lineage. Therefore, this scenario would require two HGTs as well as two independent associations and one split between the J and Fer domains, what is less parsimonious than the previous one that only requires one association and two HGT events.

Although poorly resolved as in the case of the phylogeny of the Fer domain, the phylogeny of the entire data set of J domain sequences yielded a very different picture. In fact, Viridiplantae and Thaumarchaeota did not cluster together, which was confirmed by a second analysis based on a more restricted sequence sampling. The J domains from the DnaJ-Fer proteins formed three distinct groups (indicated by colours in Figure 4B) scattered among J domain sequences of very different origins (bacterial, eukaryotic and archaeal) and being part of very diverse multidomain proteins. One group contained the J domains from Viridiplantae DnaJ-Fer proteins, another contained those from the group I.1b Thaumarchaeota (i.e., '*Ca. N. gargensis*' and *N. viennensis*), whereas group I.1a Thaumarchaeota emerged in another part of the tree (Figure 4B). This separation in three groups suggested that the J domains of the DnaJ-Fer proteins have different origins. However, this could be due just to the overall poor resolution of the trees. Thus, to discriminate between these two hypotheses (i.e. different origins or lack of phylogenetic signal) we compared the topology of the ML tree with AU tests against four constrained topologies reflecting alternative scenarios for the origin of the DnaJ domain contained in the DnaJ-Fer proteins: 1) the grouping of the J domains of the DnaJ-Fer proteins of the two groups of Thaumarchaeota I.1a and I.1b (Topology 2); 2) the monophyly of these sequences plus the J domains of the DnaJ-Fer proteins of the Viridiplantae (Topology 3); 3) the monophyly of group I.1a Thaumarchaeota and Viridiplantae DnaJ-Fer J domains

(Topology 4); and 4) the monophyly of group I.1b Thaumarchaeota and Viridiplantae DnaJ-Fer J domains (Topology 5) (Table 1), the other nodes remaining unchanged. The five topologies were used for the AU test with the alignment of J domain sequences used for the inference of the initial topology (Topology 1). All the four alternative topologies were significantly rejected ($p < 0.05$, Table 1), which indicated that the J domains found in the DnaJ-Fer proteins probably have three independent evolutionary origins.

To reconcile this observation with those from the Fer domain (see above), the most parsimonious hypothesis would be that homologous replacements of the J domain occurred twice in Thaumarchaeota after their acquisition of the DnaJ-Fer protein from Viridiplantae (Figure 4B). Such independent homologous replacements could also explain the structural differences observed between the sequences of Viridiplantae and Thaumarchaeota, namely the presence of different central regions separating the J and the Fer domains (large in group I.1a Thaumarchaeota, short in group I.1b Thaumarchaeota, and its absence in Viridiplantae, see above).

The complex evolutionary history of the DnaK/DnaJ/GrpE system in Archaea

Two of the three DnaJ-Fer proteins of *C. reinhardtii* (CDJ3 and CDJ4) have been shown to interact with the chloroplast Hsp70B proteins [24]. These proteins together with their partners, the co-chaperone DnaJ and the nucleotide exchange factor GrpE, are widely distributed in eukaryotes and also in bacteria (where they are encoded in a gene cluster). In contrast, in Archaea, they were initially reported only in lineages of mesophilic and thermophilic euryarchaeota [16,17]. However, at that time the available complete genome sequences were far from covering the whole diversity of the archaeal phyla [25] and many major lineages were not represented. Our survey of about a hundred archaeal genomes now available allowed us confirming the presence of the three genes in all members of the lineages where they were initially reported (Methanobacteriales, Thermoplasmatales, Halobacteriales and Methanosarcinales, Additional file 1) [16,17]. In addition, we also found them in many other major lineages, such as Thaumarchaeota, 'Aigarchaeota', ARMAN group, DHEV2 group, Nanohaloarchaeales, Methanomicrobiales, and Methanocellales (Additional file 1:

Table S1). This increased considerably the diversity of archaea harboring the DnaK system. Worth noting, all these archaea were either mesophilic or thermophilic organisms, underlying the absence of the DnaK/DnaJ/GrpE system in hyperthermophilic archaea. In fact, to the noticeable exception of mesophilic Methanococcales for which we identified DnaK and GrpE homologues only in *Methanococcus vannielii* SB, which were not included in our phylogenetic trees because of their extreme sequence divergence, all mesophilic and thermophilic archaea encoded these three genes and, in most of these archaeal genomes, the three genes were clustered together as occurs in Bacteria (Additional file 1).

In agreement with previous studies [13,14,16,17] our phylogenetic analysis of a subset of 136 sequences representative of the genetic diversity of bacterial, archaeal and eukaryotic DnaK/Hsp70 sequences from complete genomes supported a clear separation between eukaryotic and prokaryotic sequences (BV = 100% and PP = 1.00), and the grouping of mitochondrial and chloroplast sequences with Alphaproteobacteria (BV < 50% and PP < 0.50) and Cyanobacteria (BV = 60% and PP = 1.00), respectively (Additional file 2). By contrast, bacterial and archaeal sequences did not form two separated monophyletic groups but appeared intricately mixed, suggesting that HGT occurred between these two domains of life. To increase the resolution of the evolutionary relationships among prokaryotic DnaK proteins, we reanalysed this dataset after removing the eukaryotic sequences (Figure 5). The monophyly of most bacterial phyla was recovered, often with strong statistical support: Aquificae (BV = 100%; PP = 1.00); Cyanobacteria (BV = 81%; PP = 1.00, a second copy that groups with *Deinococcus/Thermus* exists in some Cyanobacteria); Actinobacteria (BV = 100%; PP = 1.00); Thermotogae (BV = 100%; PP = 1.00); Dictyoglomi (BV = 100%; PP = 1.00); *Deinococcus/Thermus* (BV = 98%; PP = 0.82); Spirochaetes (BV = 58%; PP = 0.68); Chlamydiae and Verrucomicrobia (BV < 50%; PP = 0.78); Alpha-, Beta-, Gamma-, Deltaproteobacteria (BV < 50%; PP = 1.00). Similarly, the monophyly of most archaeal orders and classes harbouring a DnaK gene was recovered with high support: Halobacteria and Nanohaloarchaea (BV = 100% ; PP = 1.00); Methanosarcinales (except *Methanosaeta thermophila*, BV = 100%; PP = 1.00); Methanomicrobiales (BV = 100%; PP = 1.00); Methanobacteriales (BV = 100%; PP = 1.00); Methanocellales (BV = 100%; PP = 1.00); ARMAN (BV = 52%; PP = 1.00); Thermoplasmatales (BV = 100%; PP = 1.00) together with *Aciduliprofundum boonei* (BV = 88%; PP = 1.00) and Thaumarchaeota (BV = 95%; PP = 1.00). This indicated the ancestral presence of DnaK in these groups (i.e. prior to their diversification) and that very few HGTs among them occurred after their diversification.

As in most molecular phylogenies, the relationships among bacterial phyla remained mostly unresolved

(BV < 50% and PP < 0.95, Figure 5). However, the ancestral presence of DnaK in most of them suggested that this protein was present in the last common ancestor of bacteria. By contrast, the relationships among archaeal orders and classes were well resolved but in strong contradiction with the reference species phylogeny of this domain [25]. To assess the robustness of this contradiction, we tested if the reference archaeal phylogeny was significantly rejected by the archaeal DnaK dataset. For that, the archaeal relationships observed in the DnaK ML tree were compared to those of the archaeal species reference phylogeny [25]. The AU test indicated that the DnaK dataset strongly rejected the reference phylogeny (p = 0.0). In agreement with previous proposals [16,17], this supported the hypothesis that DnaK was not present in the ancestor of Archaea and that it was acquired secondarily by some members of this domain by HGT from bacteria. A careful examination of the DnaK trees suggested that at least two independent inter-domain HGT events occurred: one to the ancestor of Halobacteria and Nanohaloarchaea and another to the ancestor of Class II methanogens (i.e., Methanosarcinales, Methanomicrobiales and Methanocellales [25]). However, the phylogeny of DnaK was not resolved enough to determine without ambiguities the bacterial donors at the origin of these HGT even if, as previously proposed [17], Firmicutes may be a possible donor in the case of the HGT to methanogenic archaea. The initial acquisitions were probably followed by secondary HGTs to and among Methanobacteriales, ARMAN, Thermoplasmatales + DHEV2, '*Ca. Caldiarchaeum subterraneum*' and Thaumarchaeota. Worth noting, the two latter lineages did not form a monophyletic group (Figure 5), in contradiction with the expected species phylogeny [25]. This suggests two independent acquisitions of the gene coding for DnaK from two different euryarchaeotal donors, but the statistical support for the corresponding branches are too low to reach a definitive conclusion. Interestingly, if inter-domain HGTs from Bacteria to Archaea were clearly supported by our analyses, at least one HGT occurred in the opposite direction. This concerned the DnaK of *Elusimicrobium minutum* (Class Elusimicrobia, previously referred as Termite Group 1), which was nested among archaeal sequences with strong support (Figure 5). This ultramicrobacterium was isolated from humivorous beetle larvae and some close relatives have been detected in gut or faeces of termites, cockroaches, and mammals such as chimpanzee, horses or cows [36], environments that are also inhabited by diverse methanogenic archaea. Accordingly, HGT among those microorganisms is not unexpected.

The phylogenies of GrpE and DnaJ were less resolved than that of DnaK, in particular for the deepest nodes (Additional file 3 and Additional file 4), as expected from the smaller number of conserved positions that

Table 1 AU tests of scenarios for the origin of the J domain contained in DnaJ-Fer proteins

Topology	Scenario	p-value
1	Viridiplantae, Group I.1a and Group I.1b not monophyletic	0.991
2	Group I.1a and Group I.1b monophyletic	0
3	Viridiplantae, Group I.1a, and Group I.1b monophyletic	0
4	Viridiplantae and Group I.1a monophyletic	0
5	Viridiplantae and Group I.1b monophyletic	0.022

could be kept for phylogenetic analyses (105 and 227 positions, respectively). However, despite the weaker signal, most of the monophyletic groups observed in the DnaK trees were also recovered in the GrpE and DnaJ phylogenies, which suggested that the three proteins have undergone similar evolutionary histories. This agreed with the fact that the corresponding genes are clustered in most prokaryotic genomes (Additional file 1).

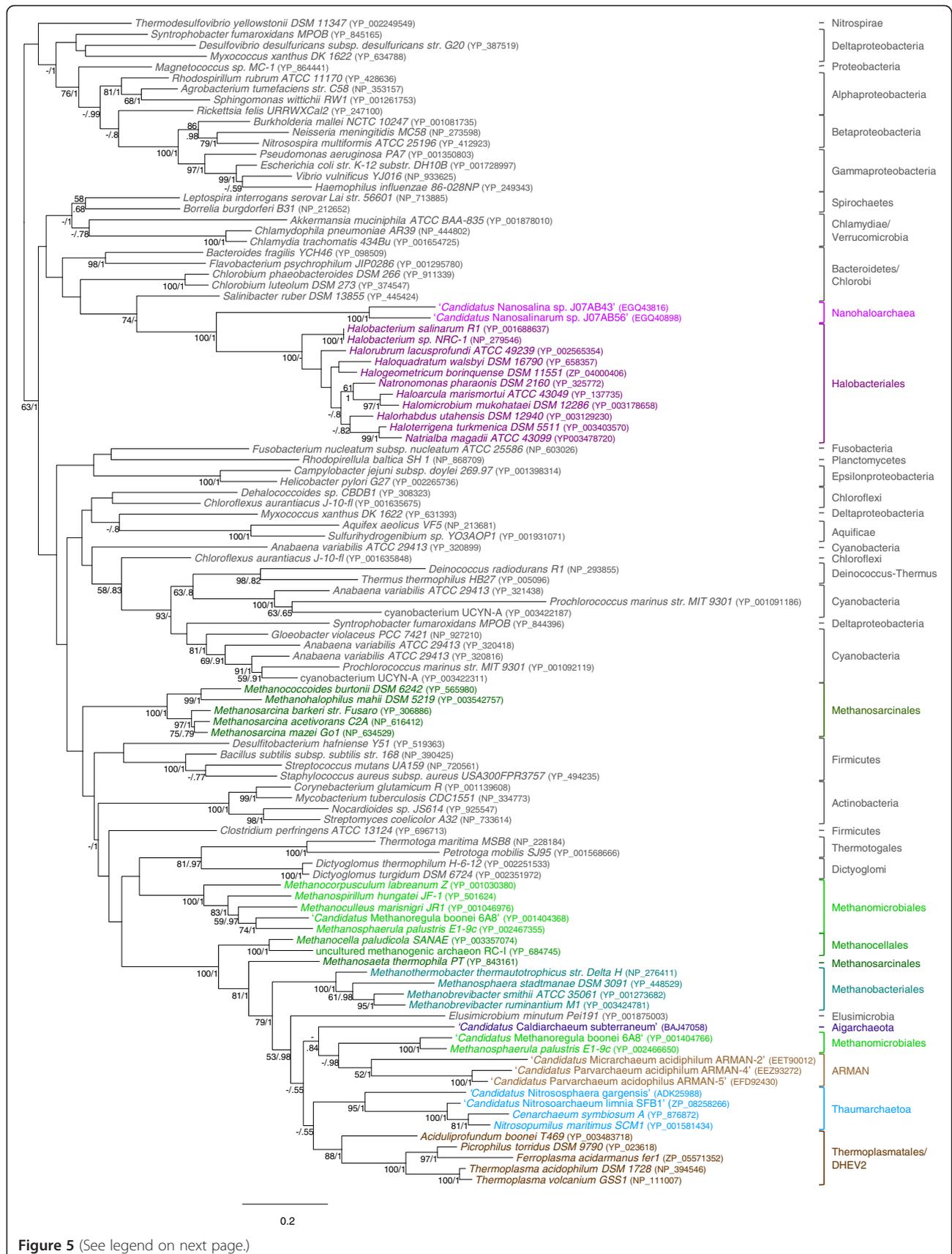
Discussion

The unexpected discovery eight years ago of an atypical protein composed of Ferredoxin domain associated to a J domain in the marine archaeal fosmid EC-39 led to the proposal that Thaumarchaeota should have a DnaK protein [18]. This prediction was confirmed a few years later after the identification of genes coding for the DnaK/DnaJ/GrpE system in complete genome sequences of representatives of this phylum [37]. Our phylogenomic analysis showed that the DnaJ-Fer and the DnaK/DnaJ/GrpE proteins have two different origins in this archaeal phylum. The thaumarchaeotal DnaJ-Fer protein resulted from a complex history involving at least two inter-domain HGTs: from Cyanobacteria to Viridiplantae and then from Viridiplantae to Thaumarchaeota, in addition to two independent replacements of the original viridiplantae J domain by J domains of unknown bacterial origin during the diversification of Thaumarchaeota (Figure 3B). By contrast, the phylogenetic analysis of the DnaK, DnaJ and GrpE proteins suggested that these proteins were acquired by the ancestor of Thaumarchaeota by HGT from an unidentified euryarchaeotal donor. The thaumarchaeotal DnaK/DnaJ/GrpE and DnaJ-Fer form therefore a complex chimera, mixing components from bacterial, euryarchaeotal and eukaryotic origin. Identifying the precise functional role of these proteins in Thaumarchaeota will require further investigation.

Due to its large distribution in present day organisms, DnaK was initially proposed as a good phylogenetic marker to infer ancient phylogenies and, more precisely, to decipher the relationships among the three domains of life and their main phyla [38-40]. In particular, the presence of a 23 amino acids deletion shared by Actinobacteria and Firmicutes (referred as Monoderma) and Archaea, but not by Gram negative bacteria (referred as Diderma) and Eucarya was interpreted as the evidence

that Archaea derived from Gram positive bacteria [38], dismissing the hypothesis of inter-domain HGT. However, the phylogenetic analysis of DnaK (and of its two partners DnaJ and GrpE) showed later that the evolutionary history of these proteins has been largely affected by HGT. In particular, the strong discrepancies observed between the phylogeny of archaeal DnaK and the species tree indicates that multiple HGTs are responsible of the taxonomic distribution of DnaK in Archaea [17,41]. Thus, the deletion detected by Gupta (which now has been shown to be present also in Thermotogae, Dictyoglomi and Fusobacteria) should not be interpreted as relevant for species phylogeny but just as a strong signature for the gene transfers mentioned above. Therefore, DnaK appears not to be a reliable marker to infer ancient evolutionary relationships, as already suggested in previous works [41].

Phylogenetic and molecular analyses have shown that adaptation to mesophily occurred several times independently during the diversification of Archaea [22,42,43]. Because, DnaK is an important heat shock chaperone involved, among others, in thermal stresses [1], we have postulated previously that its acquisition could have helped for the adaptation of archaeal hyperthermophiles to mesophilic environments [18]. Strengthening this hypothesis, DnaK, DnaJ and GrpE are found only in thermophilic and mesophilic archaea which have acquired these genes several times independently, either from bacteria through inter-domain HGTs or from archaea already adapted to mesophilic environments. However, the acquisition of DnaK cannot be considered as obligatory for that adaptation. For example, in the case of Methanococcales, we identified an atypical DnaK gene in *Methanococcus vannielii* SB, but we did not detect it in other mesophilic members of this archaeal order. This illustrated the fact that additional contributing factors, such as protein amino acid composition, are surely important to determine the optimal growth temperature of microorganisms. For instance, always among the Methanococcales, the mesophilic *Methanococcus maripaludis* has been shown to harbour amino acid signatures typical of thermophilic or hyperthermophilic organisms [44], making tempting to speculate that the proteins of this archaeon are intrinsically resistant to heat shocks and, thus, that heat shock chaperones are dispensable in this organism.



(See figure on previous page.)

Figure 5 Unrooted ML tree of the DnaK protein. The tree was inferred with TreeFinder with the LG + Γ 4 model (107 sequences and 444 positions). Numbers at nodes represent bootstrap values (100 replicates of the original alignment) and Bayesian posterior probabilities computed with TreeFinder and MrBayes, respectively (only values >50% and 0.5 are shown, dashes indicate that the corresponding support is inferior to the threshold, whereas when both supports are inferior to the thresholds no support values are indicated). The scale bar represents the average number of substitutions per site.

Beside those aspects, the evolutionary history of the DnaJ-Fer protein provided an interesting temporal landmark between the Eucarya and Archaea domains. Indeed, the association between the Fer and the J domains composing this protein very likely occurred in an ancestor of the Viridiplantae, before their diversification but after their divergence from the two other Archaeplastida lineages (i.e. the Glaucophyta and the Rhodophyta, which do not have this fused protein) (Figure 3A). Then, the resulting gene was transferred to the ancestor of Thaumarchaeota groups I.1a and I.1b (Figure 3B), more precisely before the divergence of these two lineages but likely after their separation from the HWCG III group (Figure 3A). This indicated that the divergence of the groups I.1a and I.1b is more recent than the divergence between Viridiplantae and the two other Archaeplastida lineages but more ancient than the diversification of Viridiplantae. This illustrates how HGTs can be useful to date evolutionary events relatively against each other [45]. According to fossil record and molecular dating estimates, the divergence of Viridiplantae from the two other Archaeplastida lineages occurred ~950 million years ago whereas the diversification of Viridiplantae started ~750 million years ago [46]. The HGT from Viridiplantae to Thaumarchaeota occurred most likely during this time window, so the divergence of the groups I.1a and I.1b Thaumarchaeota and their diversification could be less than ~950 million years old.

Conclusions

Phylogenomic analysis supports that the proteins DnaK, DnaJ, GrpE and DnaJ-Fer have a chimerical origin in Thaumarchaeota, which acquired them by HGT from different donors, including bacterial and eukaryotic species. Similar HGT events have occurred independently in other archaeal groups. This suggests that the acquisition of these proteins has probably played an important role in the convergent adaptation of these archaea to mesophilic and thermophilic lifestyles from their hyperthermophilic ancestors. In addition, these HGT events can be used as markers for the relative dating of the diversification of donor and acceptor groups as, for example, the Thaumarchaeota, which have received their DnaJ-Fer protein from Archaeplastida.

Methods

Dataset assembly

The DnaJ-Fer protein homologues were retrieved from the non-redundant (nr) and the environmental databases

at the NCBI (<http://www.ncbi.nlm.nih.gov>) with the BlastP program (default parameters) [47] using as seeds the DnaJ-Fer protein from the uncultured archaeon DeepAnt-EC39 fosmid (AY316120.1), and the sequences of *Chlamydomonas reinhardtii* CDJ3, 4 and 5 (XP_001700257.1, XP_001699768.1 and XP_001700843.1 respectively). The DnaJ-Fer sequence from *Nitrososphaera viennensis* was kindly provided by C. Schleper and A. Spang. To ensure the exhaustive retrieval of eukaryote sequences we queried EST and ongoing genome project databases: the JGI (<http://genome.jgi-psf.org/>) for *Ostreococcus* sp. RCC809, *Emiliania huxleyi*, *Chlorella vulgaris* and *Chlorella* sp. NC64A; the TIGR (<http://plantta.jcvi.org/>) for *Pinus taeda*; the *Cyanidioschyzon merolae* Genome Project (<http://merolae.biol.s.u-tokyo.ac.jp/>), and the *Galdieria sulphuraria* Genome Project (<http://genomics.msu.edu/galdieria/about.html>). The absence of DnaJ-Fer homologues in any archaeal or eukaryotic complete genome was verified by tBlastN searches against the corresponding nucleic acid sequence. The presence of the J and the Fer domains in the retrieved homologues was systematically verified using Pfam (Pfam profiles PF00226 and PF13459, respectively). The J and Fer domains were then analysed separately using the same strategy as previously to retrieve proteins containing these domains.

DnaK, GrpE and DnaJ homologues were retrieved from 92 archaeal complete genome sequences available at NCBI with BlastP (default parameters) using the sequences from *N. maritimus* as seeds (YP_001581434, YP_001581433.1 and YP_001581435, respectively). The absence of homologues in any genome was systematically verified by tBlastN searches against the corresponding nucleic acid sequence. Eukaryotic and bacterial homologues were retrieved from a subset of four and 86 complete genomes representative of the taxonomic diversity of these two domains using BlastP (default parameters). In the case of DnaJ, we checked the domain composition of the retrieved homologues with PFAM in order to distinguish *bona fide* DnaJ proteins (harbouring a J-domain (PF00226), the cysteine rich central domain (PF00684) and the C-terminal domain (PF01556)) from other J-proteins.

We thus obtained six different sequence datasets, and we tested various programs to align them, including (Mafft v6.833b [48], Probcons v1.12 [49], and Muscle v3.6 [50]). The quality of the resulting alignments was visually inspected in order to keep those for which the residues of the conserved domains were correctly aligned.

Probcons provided better results for the DnaJ-Fer, the J domain and the Fer domain datasets, whereas Mafft provided better results in the case of the DnaK, GrpE and DnaJ datasets. The selected alignments were edited and manually refined with the program ED of the MUST package [51]. The regions where the alignment was ambiguous were removed using the NET program from the MUST package.

Phylogenetic reconstruction

The DnaJ-Fer, DnaK, GrpE and DnaJ alignments were analysed by maximum likelihood (ML) and Bayesian approaches (BI). For each dataset, the LG model was proposed as the best suited evolutionary model according to the "propose model tool" of TreeFinder v2011 [52] with the AICc criterion. Alternative models (e.g. WAG, JTT, etc.) were also tested. The resulting trees were consistent with those inferred with the LG model (not shown).

ML tree reconstructions were performed using PhyML v3.0 [53] and TreeFinder v2011 [52]. The robustness of the resulting trees was estimated by the non-parametric procedure implemented in PhyML and TreeFinder (100 replicates of the original dataset). The resulting trees were very similar, so we decided to show only the ML trees inferred with TreeFinder.

BI of DnaJ-Fer, DnaK, DnaJ and GrpE proteins was carried out with PhyloBayes v 3.3 with the LG model and a gamma distribution of substitution rates with four categories [54]. Phylobayes was run with four independent chains for at least 10,000 cycles, saving one tree in ten. The first 300 trees were discarded as "burn-in", and the remaining trees from each chain were used to test for convergence and compute the 50% majority rule consensus tree. In the case of DnaK, the chains did not converge even after 10,000 cycles. Therefore, BI trees for this marker were computed with MrBayes v.3.0b4 [55] with a mixed substitution model and a Gamma distribution of substitution rates with 4 categories. Searches were run with 4 chains of 1,000,000 generations for which the first 1,000 generations were discarded as "burn in", trees being sampled every 100 generations.

The analyses of the J and Fer domains were divided in two steps. First, all the homologous sequences were analysed by neighbor-joining (NJ) using the MUST package [51]. Based on this preliminary phylogenetic tree, we selected the closest homologues of the DeepAnt-EC39 fosmid and *C. reinhardtii* sequences and a subset of more distantly related homologues representative of the genetic diversity of these domains. These sequences were used to carry out ML and BI analysis with TreeFinder, PhyML and PhyloBayes as previously described.

The comparison of different tree topologies was done by applying the Approximately Unbiased test [56] implemented

in TreeFinder with the same evolutionary models and parameters as for ML phylogenetic inference.

Additional files

Additional file 1: Table showing the taxonomic distribution of DnaJ-Fer, DnaK, DnaJ and GrpE proteins in Archaea. Numbers correspond to accession numbers in the NCBI Genpep database in the 92 complete genomes available in July 2011 and that of, '*C. Nitrosoarchaeum koreensis*', a thaumarchaeotal genome available more recently. The two divergent sequences of DnaK and GrpE found in *Methanococcus vannielii* SB are underlined.

Additional file 2: Unrooted ML tree of the DnaK protein (136 sequences and 444 positions) inferred with TreeFinder and the LG + Γ 4 model. Numbers at nodes represent bootstrap values and Bayesian posterior probabilities computed with TreeFinder and MrBayes, respectively (only values >50% and 0.5 are shown, dashes indicate the corresponding support is inferior to the threshold, whereas when both supports are inferior to the thresholds no support values are indicated). Archaeal sequences are shown with colours according to their taxonomic classification. The scale bar represents the average number of substitutions per site.

Additional file 3: Unrooted ML tree of the DnaJ protein (102 sequences and 227 positions) inferred with TreeFinder and the LG + Γ 4. Numbers at nodes represent bootstrap values and Bayesian posterior probabilities computed with TreeFinder and MrBayes, respectively (only values >50% and 0.5 are shown, dashes indicate that the corresponding support is inferior to the threshold, whereas when both supports are inferior to the thresholds no support values are indicated). Archaeal sequences are shown with colours according to their taxonomic classification. The scale bar represents the average number of substitutions per site.

Additional file 4: Unrooted ML tree of the GrpE protein (101 sequences and 105 positions) inferred with TreeFinder and the LG + Γ 4. Numbers at nodes represent bootstrap values and Bayesian posterior probabilities computed with TreeFinder and MrBayes, respectively (only values >50% and 0.5 are shown, dashes indicate the corresponding value is inferior to the threshold, whereas when both supports are inferior to the thresholds no support values are indicated). Archaeal sequences are shown with colours according to their taxonomic classification. The scale bar represents the average number of substitutions per site.

Competing interests

The authors have declared that no competing interests exist.

Authors' contributions

CB-A and DM conceived this study. CP, CB-A and DM designed and carried out the phylogenetic analyses. CP, CB-A, PL-G and DM wrote the manuscript. All authors read and approved the final manuscript.

Acknowledgements

This work was supported by the Agence Nationale de la Recherche (ANR EvolDeep; contract number ANR-08-GENM-024-002). C. Petitjean was the recipient of a grant from the ANR EvolDeep. C. Brochier-Armanet was member of the Institut Universitaire de France, and was funded by an ATIP from the Centre National de la Recherche Scientifique (CNRS) and by the Ancestrôme project (ANR-10-BINF-01-01) and by the ANR-10-BINF-01-01 "Ancestrôme" grant. We acknowledge C. Schleper, A. Spang and J. de la Torre for sharing unpublished data.

Author details

¹UPR CNRS 9043, Laboratoire de Chimie Bactérienne, Université d'Aix-Marseille (AMU), 13402 Marseille, Cedex 20, France. ²UMR CNRS 8079, Institut d'Ecologie, Systématique et Evolution Université Paris-Sud, 91405 Orsay, Cedex, France. ³CNRS, UMR5558, Laboratoire de Biométrie et Biologie Evolutive, Université de Lyon, Université Lyon 1, 43 boulevard du 11 novembre 1918, 69622, Villeurbanne, France.

Received: 19 March 2012 Accepted: 25 October 2012
Published: 26 November 2012

References

- Mayer MP, Bukau B: **Hsp70 chaperones: cellular functions and molecular mechanism.** *Cell Mol Life Sci* 2005, **62**(6):670–684.
- Young JC: **Mechanisms of the Hsp70 chaperone system.** *Biochem Cell Biol* 2010, **88**(2):291–300.
- Morano KA: **New tricks for an old dog: the evolving world of Hsp70.** *Ann N Y Acad Sci* 2007, **1113**:1–14.
- Kampinga HH, Craig EA: **The HSP70 chaperone machinery: J proteins as drivers of functional specificity.** *Nat Rev Mol Cell Biol* 2010, **11**(8):579–592.
- Harrison C: **GrpE, a nucleotide exchange factor for DnaK.** *Cell Stress Chaperones* 2003, **8**(3):218–224.
- Laloraya S, Gambill BD, Craig EA: **A role for a eukaryotic GrpE-related protein, Mge1p, in protein translocation.** *Proc Natl Acad Sci USA* 1994, **91**(14):6481–6485.
- Schroda M, Vallon O, Whitelegge JP, Beck CF, Wollman FA: **The chloroplastic GrpE homolog of Chlamydomonas: two isoforms generated by differential splicing.** *Plant Cell* 2001, **13**(12):2823–2839.
- Alberti S, Esser C, Hohfeld J: **BAG-1—a nucleotide exchange factor of Hsc70 with multiple cellular functions.** *Cell Stress Chaperones* 2003, **8**(3):225–231.
- Craig EA, Huang P, Aron R, Andrew A: **The diverse roles of J-proteins, the obligate Hsp70 co-chaperone.** *Rev Physiol Biochem Pharmacol* 2006, **156**:1–21.
- Liberek K, Marszałek J, Ang D, Georgopoulos C, Zyllicz M: **Escherichia coli DnaJ and GrpE heat shock proteins jointly stimulate ATPase activity of DnaK.** *Proc Natl Acad Sci USA* 1991, **88**(7):2874–2878.
- Walsh P, Bursac D, Law YC, Cyr D, Lithgow T: **The J-protein family: modulating protein assembly, disassembly and translocation.** *EMBO Rep* 2004, **5**(6):567–571.
- Qiu XB, Shao YM, Miao S, Wang L: **The diversity of the DnaJ/Hsp40 family, the crucial partners for Hsp70 chaperones.** *Cell Mol Life Sci* 2006, **63**(22):2560–2570.
- Boorstein WR, Ziegelhoffer T, Craig EA: **Molecular evolution of the HSP70 multigene family.** *J Mol Evol* 1994, **38**(1):1–17.
- Renner T, Waters ER: **Comparative genomic analysis of the Hsp70s from five diverse photosynthetic eukaryotes.** *Cell Stress Chaperones* 2007, **12**(2):172–185.
- Nordhues A, Miller SM, Muhlhaus T, Schroda M: **New insights into the roles of molecular chaperones in Chlamydomonas and Volvox.** *Int Rev Cell Mol Biol* 2010, **285**:75–113.
- Macario AJ, Brocchieri L, Shenoy AR, Conway de Macario E: **Evolution of a protein-folding machine: genomic and evolutionary analyses reveal three lineages of the archaeal hsp70(dnaK) gene.** *J Mol Evol* 2006, **63**(1):74–86.
- Gribaldo S, Lumia V, Creti R, de Macario EC, Sanangelantoni A, Cammarano P: **Discontinuous occurrence of the hsp70 (dnaK) gene among Archaea and sequence features of HSP70 suggest a novel outlook on phylogenies inferred from this protein.** *J Bacteriol* 1999, **181**(2):434–443.
- Lopez-Garcia P, Brochier C, Moreira D, Rodriguez-Valera F: **Comparative analysis of a genome fragment of an uncultivated mesopelagic crenarchaeote reveals multiple horizontal gene transfers.** *Environ Microbiol* 2004, **6**(1):19–34.
- DeLong EF: **Archaea in coastal marine environments.** *Proc Natl Acad Sci USA* 1992, **89**(12):5685–5689.
- Fuhrman JA, McCallum K, Davis AA: **Novel major archaeobacterial group from marine plankton.** *Nature* 1992, **356**(6365):148–149.
- Brochier-Armanet C, Boussau B, Gribaldo S, Forterre P: **Mesophilic Crenarchaeota: proposal for a third archaeal phylum, the Thaumarchaeota.** *Nat Rev Microbiol* 2008, **6**(3):245–252.
- Brochier-Armanet C, Gribaldo S, Forterre P: **Spotlight on the Thaumarchaeota.** *ISME J* 2012, **6**(2):227–230.
- Pester M, Schleper C, Wagner M: **The Thaumarchaeota: an emerging view of their phylogeny and ecophysiology.** *Curr Opin Microbiol* 2011, **14**(3):300–306.
- Dorn KV, Willmund F, Schwarz C, Henselmann C, Pohl T, Hess B, Veyel D, Usadel B, Friedrich T, Nickelsen J, et al: **Chloroplast DnaJ-like proteins 3 and 4 (CDJ3/4) from Chlamydomonas reinhardtii contain redox-active Fe-S clusters and interact with stromal HSP70B.** *Biochem J* 2010, **427**(2):205–215.
- Brochier-Armanet C, Forterre P, Gribaldo S: **Phylogeny and evolution of the Archaea: one hundred genomes later.** *Curr Opin Microbiol* 2011, **14**(3):274–281.
- Adl SM, Simpson AG, Farmer MA, Andersen RA, Anderson OR, Barta JR, Bowser SS, Brugerolle G, Fensome RA, Fredericq S, et al: **The new higher level classification of eukaryotes with emphasis on the taxonomy of protists.** *J Eukaryot Microbiol* 2005, **52**(5):399–451.
- Hallam SJ, Konstantinidis KT, Putnam N, Schleper C, Watanabe Y, Sugahara J, Preston C, de la Torre J, Richardson PM, DeLong EF: **Genomic analysis of the uncultivated marine crenarchaeote Cenarchaeum symbiosum.** *Proc Natl Acad Sci USA* 2006, **103**(48):18296–18301.
- Walker CB, de la Torre JR, Klotz MG, Urakawa H, Pinel N, Arp DJ, Brochier-Armanet C, Chain PS, Chan PP, Gollabgir A, et al: **Nitrosopumilus maritimus genome reveals unique mechanisms for nitrification and autotrophy in globally distributed marine crenarchaea.** *Proc Natl Acad Sci USA* 2010, **107**(19):8818–8823.
- Blainey PC, Mosier AC, Potanina A, Francis CA, Quake SR: **Genome of a Low-Salinity Ammonia-Oxidizing Archaeon Determined by Single-Cell and Metagenomic Analysis.** *PLoS One* 2011, **6**(2):e16626.
- Kim BK, Jung MY, Yu DS, Park SJ, Oh TK, Rhee SK, Kim JF: **Genome sequence of an ammonia-oxidizing soil archaeon, "Candidatus Nitrosoarchaeum koreensis" MY1.** *J Bacteriol* 2011, **193**(19):5539–5540.
- Spang A, Poehlein A, Offire P, Zumbagel S, Haider S, Rychlik N, Nowka B, Schmeisser C, Lebedeva EV, Rattai T, et al: **The genome of the ammonia-oxidizing Candidatus Nitrososphaera gargensis: insights into metabolic versatility and environmental adaptations.** *Environmental microbiology* 2012.
- Nunoura T, Takaki Y, Kakuta J, Nishi S, Sugahara J, Kazama H, Chee GJ, Hattori M, Kanai A, Atomi H, et al: **Insights into the evolution of Archaea and eukaryotic protein modifier systems revealed by the genome of a novel archaeal group.** *Nucleic Acids Res* 2011, **39**(8):3204–3223.
- Keeling PJ: **The endosymbiotic origin, diversification and fate of plastids.** *Philos Trans R Soc Lond B Biol Sci* 2010, **365**(1541):729–748.
- Deschamps P, Moreira D: **Signal conflicts in the phylogeny of the primary photosynthetic eukaryotes.** *Mol Biol Evol* 2009, **26**(12):2745–2753.
- Criscuolo A, Gribaldo S: **Large-scale phylogenomic analyses indicate a deep origin of primary plastids within cyanobacteria.** *Mol Biol Evol* 2011, **28**(11):3019–3032.
- Geissinger O, Herlemann DP, Morschel E, Maier UG, Brune A: **The ultramicrobacterium "Elusimicrobium minutum" gen. nov., sp. nov., the first cultivated representative of the termite group 1 phylum.** *Appl Environ Microbiol* 2009, **75**(9):2831–2840.
- Spang A, Hatzepichler R, Brochier-Armanet C, Rattai T, Tischler P, Spieck E, Streit W, Stahl DA, Wagner M, Schleper C: **Distinct gene set in two different lineages of ammonia-oxidizing archaea supports the phylum Thaumarchaeota.** *Trends Microbiol* 2010, **18**(8):331–340.
- Gupta RS: **What are archaeobacteria: life's third domain or monoderm prokaryotes related to Gram-positive bacteria? A new proposal for the classification of prokaryotic organisms.** *Mol Microbiol* 1998, **229**(3):695–708.
- Griffiths E, Gupta RS: **The use of signature sequences in different proteins to determine the relative branching order of bacterial divisions: evidence that Fibrobacter diverged at a similar time to Chlamydia and the Cytophaga-Flavobacterium-Bacteroides division.** *Microbiology* 2001, **147**(Pt 9):2611–2622.
- Gupta RS: **Origin of diderm (Gram-negative) bacteria: antibiotic selection pressure rather than endosymbiosis likely led to the evolution of bacterial cells with two membranes.** *Antonie Van Leeuwenhoek* 2011, **100**(2):171–182.
- Philippe H, Budin K, Moreira D: **Horizontal transfers confuse the prokaryotic phylogeny based on the HSP70 protein family.** *Mol Microbiol* 1999, **31**(3):1007–1009.
- Gribaldo S, Brochier-Armanet C: **The origin and evolution of Archaea: a state of the art.** *Philos Trans R Soc Lond B Biol Sci* 2006, **361**(1470):1007–1022.
- Groussin M, Gouy M: **Adaptation to environmental temperature is a major determinant of molecular evolutionary rates in archaea.** *Mol Biol Evol* 2011, **28**(9):2661–2674.
- Puigbo P, Pasamontes A, Garcia-Valve S: **Gaining and losing the thermophilic adaptation in prokaryotes.** *Trends in genetics: TIG* 2008, **24**(1):10–14.
- Huang J, Gogarten JP: **Ancient horizontal gene transfer can benefit phylogenetic reconstruction.** *Trends in genetics: TIG* 2006, **22**(7):361–366.
- Douzery EJ, Snell EA, Baptiste E, Delsuc F, Philippe H: **The timing of eukaryotic evolution: does a relaxed molecular clock reconcile proteins and fossils?** *Proc Natl Acad Sci USA* 2004, **101**(43):15386–15391.
- Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ: **Gapped BLAST and PSI-BLAST: a new generation of protein database search programs.** *Nucleic Acids Res* 1997, **25**(17):3389–3402.
- Katoh K, Misawa K, Kuma K, Miyata T: **MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform.** *Nucleic Acids Res* 2002, **30**(14):3059–3066.

49. Do CB, Mahabhashyam MS, Brudno M, Batzoglou S: **ProbCons: Probabilistic consistency-based multiple sequence alignment.** *Genome Res* 2005, **15**(2):330–340.
50. Edgar RC: **MUSCLE: multiple sequence alignment with high accuracy and high throughput.** *Nucleic Acids Res* 2004, **32**(5):1792–1797.
51. Philippe H: **MUST, a computer package of Management Utilities for Sequences and Trees.** *Nucleic Acids Res* 1993, **21**(22):5264–5272.
52. Jobb G, von Haeseler A, Strimmer K: **TREEFINDER: a powerful graphical analysis environment for molecular phylogenetics.** *BMC Evol Biol* 2004, **4**:18.
53. Guindon S, Gascuel O: **A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood.** *Syst Biol* 2003, **52**(5):696–704.
54. Lartillot N, Lepage T, Blanquart S: **PhyloBayes 3: a Bayesian software package for phylogenetic reconstruction and molecular dating.** *Bioinformatics* 2009, **25**(17):2286–2288.
55. Ronquist F, Huelsenbeck JP: **MrBayes 3: Bayesian phylogenetic inference under mixed models.** *Bioinformatics* 2003, **19**(12):1572–1574.
56. Shimodaira H: **An approximately unbiased test of phylogenetic tree selection.** *Syst Biol* 2002, **51**(3):492–508.

doi:10.1186/1471-2148-12-226

Cite this article as: Petitjean *et al.*: Horizontal gene transfer of a chloroplast DnaJ-Fer protein to Thaumarchaeota and the evolutionary history of the DnaK chaperone system in Archaea. *BMC Evolutionary Biology* 2012 **12**:226.

Submit your next manuscript to BioMed Central and take full advantage of:

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at
www.biomedcentral.com/submit



C. Éléments de discussion

La masse de données générées par notre analyse de l'ensemble des protéines codées dans les génomes de *C. symbiosum*, *N. maritimus* et '*Ca. Caldiarchaeum subterraneum*' n'a pas encore été étudiée et, de par son ampleur, pourrait faire l'objet d'un travail à part entière. Le cas de la protéine DnaJ-Fer n'est qu'un premier essai d'exploitation de cette information. Parmi les 3085 protéines dont la phylogénie a pu être inférée et observée, j'ai extrait quelques observations qui pourraient être exploitées. Néanmoins, je voudrais insister sur le fait que, d'une part, les jeux de données construits pour chaque protéine contiennent au plus les 200 séquences les plus proches retrouvées par BLAST. Dans certains cas, la présence de certains homologues plus lointains (particulièrement bactériens ou eucaryotes) n'a donc pu être détectée. D'autre part, l'analyse de ces données avait pour but premier la recherche de nouveaux marqueurs pour la phylogénie des archées et non une analyse rigoureuse de génomique comparative d'archées. Tout élément ayant retenu mon attention sur une phylogénie a été noté, et ce de la façon la plus systématique et standardisée possible. Enfin, un premier tri a été effectué entre les génomes de *C. symbiosum* et *N. maritimus* pour éviter une trop grande redondance, mais, par précaution et pour ne rater aucune protéine intéressante, nous avons préféré garder le génome de '*Ca. Caldiarchaeum subterraneum*' dans son ensemble quitte à analyser la phylogénie de la même famille protéique plusieurs fois, les jeux de données redondants ayant été supprimés dans la deuxième étape de sélection des marqueurs. Il existe donc une redondance entre certains jeux de données analysés (mais qui a été éliminée pour les analyses visant à étudier la phylogénie des archées et la place de la racine). Les chiffres donnés sont donc à considérer avec précaution et doivent être pris comme un ordre de grandeur et non pas comme un chiffre définitif.

On trouve ainsi un peu moins de 800 protéines pour lesquelles aucun homologue bactérien ou eucaryote n'était présent, dont 135 spécifiques aux thaumarchées et 16 spécifiques à '*Ca. Caldiarchaeum subterraneum*'. Ces protéines spécifiques aux archées au premier abord pourraient faire l'objet d'une étude et apporter de nombreux éléments quant à l'évolution du domaine Archaea et aux innovations évolutives apparues spécifiquement dans cette lignée. Parallèlement, 650 protéines sont présentes chez au moins un représentant de chaque phylum archée (Euryarchaeota, Crenarchaeota, Aigarchaeota, Thaumarchaeota, et Korarchaeota). Nous avons aussi observé la présence d'euryarchées et de crenarchées dans environ deux tiers des jeux de données, la présence de '*Ca. Korarchaeum cryptophilum*', d'une part, et des thaumarchées et de l'aigarchée, d'autre part, dans environ un tiers de ces jeux de données. Ces protéines partagées entre archées peuvent être la cible d'une étude du contenu en gènes et de la nature du dernier ancêtre commun aux archées (« LACA » pour « Last Archaeal Common Ancestor »). En effet, la présence d'une protéine dans au

Chapitre 3 : Répartition taxonomique et histoire évolutive des protéines d'archées : exemple du système chaperonne DnaK et de la protéine DnaJ-Fer.

moins deux phyla, dont notamment chez les euryarchées et au moins un des autres phyla (cf. Chapitre 2), pourrait être due à sa présence chez LACA.

De la même manière on retrouve environ un tiers des jeux de données dans lesquels des homologues bactériens et eucaryotes sont présents. Ces protéines pourraient être des cibles privilégiées pour une étude sur le dernier ancêtre commun de tous les organismes et sur les caractères conservés depuis cet ancêtre chez les archées. En ce qui concerne particulièrement les relations entre les archées et les deux autres domaines du vivant, environ 300 jeux de données contiennent des eucaryotes mais pas de bactéries ; en parallèle, environ un autre tiers des jeux de données contiennent des homologues bactériens sans eucaryotes. L'analyse des phylogénies de ces protéines pourrait être un point de départ pour analyser les relations entre les archées et les deux autres domaines du vivant, ensemble ou indépendamment.

Il est important de souligner que la présence d'un homologue bactérien, eucaryote ou d'un autre phylum d'archées n'est pas forcément le signe de la présence de cette protéine dans l'ancêtre commun de ces deux lignées, mais peut être aussi le résultat d'un transfert horizontal. Lors de ma première analyse, j'ai noté tous les cas potentiels de transfert mais il serait nécessaire de mener une nouvelle étude de ces jeux de données pour les valider. Seul un petit nombre de transferts entre différentes lignées d'archées a été considéré dans cette première analyse car il est souvent très difficile de faire la différence entre un transfert horizontal entre archées, un cas de paralogie cachée ou un manque de résolution dans une phylogénie générée automatiquement et sans analyse poussée. C'est pourquoi je préconiserais, pour étudier ces transferts, de réanalyser les jeux de données dans lesquels plusieurs phyla d'archées sont présents mais qui n'ont pas été conservés pour l'analyse de la phylogénie des archées. Une nouvelle analyse, avec un échantillonnage taxonomique différent centré sur les archées et avec toutes les informations apportées dans ma thèse concernant la phylogénie des archées et la position de la racine de ce domaine, pourrait permettre de faire un tri plus efficace et de détecter réellement les transferts entre archées. Un problème similaire se pose en ce qui concerne les transferts horizontaux entre eucaryotes et archées. Le débat existant autour de la relation entre archées et eucaryotes et les hypothèses proposant que les eucaryotes soient en réalité issus de lignées d'archées incitent à prendre avec beaucoup de précaution les observations de transferts horizontaux potentiels entre ces deux domaines. La présence d'homologues bactériens dans la plupart des phylogénies avec des homologues eucaryotes (environ 1000 sur 1300) peut rendre plus difficile l'inférence de la position des eucaryotes de par leur importante distance évolutive (i.e. générant de longues branches) et en diminuant probablement la résolution de ces phylogénies. De même que pour les relations entre phyla d'archées, une analyse centrée sur les relations entre eucaryotes et archées avec un échantillonnage taxonomique adapté serait nécessaire

pour conclure quant aux cas de transferts entre archées et eucaryotes. Les mêmes risques existent en ce qui concerne les transferts horizontaux entre archées et bactéries et les mêmes recommandations sont donc à préconiser pour leur analyse, bien que la polémique en ce qui concerne les relations entre bactéries et archées soit moins importante actuellement. Dans mon analyse initiale, j'ai noté plus de 350 cas de transferts potentiels entre archées et bactéries (le sens de ces transferts devant être analysé en détails) dont un grand nombre impliquant les thaumarchées. Il a été proposé que l'adaptation des thaumarchées mésophiles et marines à leur environnement ait impliqué de nombreux transferts depuis des bactéries vivant dans ces même milieux (López-garcía et al. 2004; Brochier-Armanet et al. 2011), l'analyse de ces phylogénies pourrait apporter de nouvelles informations quant à ces transferts et leur rôle dans cette adaptation.

La liste de l'ensemble des protéines de départ ainsi que quelques informations sur la répartition taxonomique des homologues présents dans les jeux de données résultants et sur l'analyse des phylogénies sont disponibles dans l'Annexe 2.

Discussion

Mon travail de thèse porte sur la phylogénie des archées avec des approches de phylogénomique. Il a permis de mettre en évidence un nombre relativement important de nouveaux marqueurs utilisables pour reconstruire la phylogénie des archées et améliorer sa résolution. Dans cette discussion, j'aimerais aborder certains points qui me semblent importants dans le travail autour de la phylogénomique, particulièrement lorsque celle-ci est appliquée au monde microbien.

De l'importance de la phylogénie en biologie.

La très célèbre phrase de Dobzhansky « *Nothing in biology makes sense except in the light of evolution* » (Dobzhansky 1973)¹ est une introduction parfaite à l'importance de la phylogénie en biologie. J'ai déjà parlé de l'importance de l'étude de la phylogénie des archées pour la compréhension de leur histoire évolutive ainsi que des relations qui les lient aux deux autres domaines du vivant. D'un point de vue géologique aussi, la connaissance de l'histoire évolutive des organismes est importante, par exemple l'apparition de la méthanogenèse chez les archées qui a pu modifier les conditions atmosphériques. La reconstruction de la phylogénie des organismes est aussi extrêmement importante pour comprendre les mécanismes évolutifs sous-tendant leur évolution. En effet, l'utilisation d'une phylogénie des organismes comme référence est essentielle pour comprendre l'histoire de leurs caractères, gènes ou systèmes. Par exemple, le Chapitre 3 de ma thèse présente l'histoire évolutive de quatre protéines qui ont subi de nombreux transferts horizontaux entre les trois domaines du vivant. Les protéines DnaJ et DnaJ-Fer font partie d'une famille protéique extrêmement large composée de protéines contenant un domaine protéique dnaJ. Cette famille est le résultat de multiples événements de duplications et de pertes de gènes, mais aussi de fusion avec d'autres gènes, comme dans le cas de la protéine DnaJ-Fer. Ces différents événements évolutifs particuliers, et rares pour certains, n'ont pu être mis en évidence que par la comparaison de la phylogénie de ces protéines individuelles avec une phylogénie de référence, inférée à partir de multiples marqueurs, telle que celle que nous avons construite dans le Chapitre 1. L'on pourra objecter que cette phylogénie des organismes dite « de référence », n'est pas totalement

¹ Si cet article de Dobzhansky est très important, pour l'éloquence de son titre mais aussi pour son explication de ce qu'est l'évolution biologique et de son message clairement opposé à l'antiévolutionnisme, je pense que sa lecture attentive est importante. Il me semble empreint d'une certaine forme de gradisme, mettant l'homme « à l'apex de l'évolution », et d'une volonté de réconcilier science et religion alors que ces deux domaines n'ont ni à être opposés ni à être réconciliés, puisque ne traitant pas de la même chose. Je tenais à souligner ce point à propos d'un article cité si souvent en biologie (y compris ici).

résolue et l'existence même d'une telle phylogénie est remise en question très souvent (E Bapteste et al. 2005; Dagan and Martin 2006). La découverte de nouvelles espèces et le séquençage de nouveaux génomes viennent forcément modifier la phylogénie, et de nouvelles méthodes permettent de corriger certains biais ; en revanche nous avons montré dans le Chapitre 1 que si certains points sont toujours discutés, le signal global de la phylogénie des archées reste relativement stable depuis de nombreuses années. La phylogénie globale proposée pourra donc être considérée comme une phylogénie de référence des organismes, jusqu'à la prochaine amélioration de celle-ci.

De la sélection des données.

Mes travaux de recherche montrent qu'un grand nombre de nouveaux marqueurs phylogénétiques ont pu être identifiés pour inférer la phylogénie des archées. Ainsi nous avons presque triplé le nombre de marqueurs utilisables pour aborder cette question. Même si les 200 protéines utilisées ici ne sont pas présentes chez toutes les archées, elles contiennent un signal phylogénétique suffisant pour reconstruire l'histoire évolutive de ce domaine. Les supermatrices construites à partir de ces 200 protéines, et avec la totalité des 273 marqueurs en incluant les protéines ribosomiques, l'ARN polymérase et les facteurs de transcription conservés à l'échelle des archées, permettent de clarifier certaines relations en débat et d'apporter un meilleur soutien général à la phylogénie de ce domaine. Ces nouvelles protéines peuvent avoir des répartitions taxonomiques très différentes et elles ne sont pas forcément présentes dans la totalité des archées. Elles ont été sélectionnées sur la base de leur présence dans au moins dix espèces d'archées représentant deux phyla différents ainsi que sur la base de la monophylie des séquences d'archées et des différents ordres d'archées. La présence de ces protéines chez des bactéries ou chez des eucaryotes n'a eu aucun impact sur leur choix comme marqueurs pour la phylogénie des archées. En effet, cet ensemble de protéines a été sélectionné exclusivement pour aborder la question de la phylogénie des archées et leur analyse phylogénétique a conduit à des résultats très concluants.

Les données collectées m'ont ensuite permis d'aborder la question de la position de la racine de l'arbre des archées. Compte-tenu des possibles liens qui unissent eucaryotes et archées, je me suis orientée vers l'utilisation d'un groupe extérieur bactérien. Comme dans le cas précédent, nous avons cherché le plus grand nombre possible de protéines pouvant répondre spécifiquement à cette question. L'analyse minutieuse des phylogénies individuelles des 200 marqueurs identifiés précédemment m'a permis de mettre en évidence que 38 d'entre eux présentaient les qualités requises pour étudier la position de la racine chez les archées auxquels il a été possible de rajouter 32 protéines ribosomiques. Ici encore, une sélection de marqueurs adaptés a été nécessaire pour aborder une nouvelle problématique précise et obtenir des résultats robustes.

Pour ces deux études, nous avons sélectionné des ensembles de marqueurs différents pour répondre à des questions différentes. De plus, cette sélection s'est faite uniquement sur la base de critères phylogénétiques, sans prendre en compte a priori la fonction des gènes. En ce sens ma démarche est très différente des études précédentes, qui étaient centrées sur l'analyse d'un petit nombre de protéines informationnelles, principalement les protéines ribosomiques. Certains travaux présentés dans l'Introduction, dont le but était d'inférer l'arbre du vivant et de résoudre les relations entre les trois domaines du vivant, sont passés par des recherches automatiques de protéines présentes dans un maximum de génomes. Les gènes ainsi sélectionnés correspondaient aux fonctions les plus conservées et convergeaient en général vers une liste très réduite et centrée autour des protéines ribosomiques (Tableau 2) (Ciccarelli et al. 2006; Yutin et al. 2008). Cependant, nous avons confirmé qu'une répartition taxonomique exhaustive pour chaque protéine n'est pas nécessaire dans une approche phylogénomique car les données manquantes dans certains jeux de données peuvent être compensées par la concaténation de l'ensemble dans une supermatrice. Ainsi, l'accès à de grandes quantités de données génomiques nous permet aujourd'hui d'aller chercher d'autres protéines portant un signal intéressant, et surtout de définir un nouvel ensemble de marqueurs adaptés à chaque problématique traitée.

Par exemple, les relations entre eucaryotes et archées sont très débattues actuellement (Simonetta Gribaldo et al. 2010). Le placement de la racine de l'arbre des archées en utilisant les bactéries comme groupe extérieur nous a permis d'obtenir une phylogénie racinée des archées robuste. Considérant cette phylogénie racinée comme point de départ, il serait maintenant envisageable de chercher un nouvel ensemble de marqueurs adaptés à cette question précise. Les protéines sélectionnées n'auraient pas besoin d'être présentes chez les bactéries et pourraient être choisies avec la même méthodologie que celle utilisée dans le Chapitre 2. Cette idée a déjà été proposée par Gribaldo et collaborateurs en 2010 (Simonetta Gribaldo et al. 2010), et pourrait être mise en oeuvre.

De l'augmentation de la quantité de données génomiques.

Le nombre de génomes d'archées publiés est en augmentation exponentielle depuis le début des années 2000 et les dernières avancées en termes de séquençage de cellules uniques commencent à faire exploser les quantités de données disponibles (Rinke et al. 2013). Il m'est ainsi arrivé deux fois d'entendre au cours d'une conversation, de façon plus ou moins sérieuse, que la quantité de données (génomiques ou scientifiques de façon générale) était suffisante actuellement pour répondre à toutes les questions et qu'il n'était pas forcément nécessaire d'en produire plus. Je trouve plutôt dommage d'avoir ce point de vue. Certes, l'augmentation du nombre de génomes disponibles

est complexe à gérer pour de nombreuses raisons, mais c'est aussi une mine d'information exploitable du plus grand intérêt. Il est possible que seule une petite part de ces données soit utilisée un jour, mais la génération de la totalité est très utile pour pouvoir trouver des informations intéressantes pour répondre à de nouvelles questions. Il ne faut pas oublier que de nombreux champs de la biologie sont à même d'utiliser ces séquences pour des études tout à fait différentes. Par exemple, une véritable génomique des populations commence à devenir possible pour les microorganismes, mais aussi des études très détaillées sur des caractères métaboliques, écologiques, fonctionnels ou évolutifs. Dans mon cas, la sélection des 200 nouveaux marqueurs pour la phylogénie des archées, puis des 38 pour la position de la racine, n'a pu être faite que parce qu'un nombre suffisant de génomes complets était disponible pour les archées mais aussi pour les bactéries et les eucaryotes. Avoir une répartition taxonomique représentative de la diversité de chaque domaine était nécessaire pour révéler les cas de transferts horizontaux ou de paralogies. Et dans la mesure où la recherche de ces marqueurs était par définition sans a priori, il aurait été impossible de les cibler pour un séquençage individuel. L'accès aux séquences génomiques complètes était donc nécessaire.

Parallèlement, lors de mes travaux de recherche j'ai été confrontée de nombreuses fois à l'artefact d'attraction de longues branches dû à l'évolution rapide de certains génomes. L'augmentation du nombre de génomes disponibles permet un choix plus intéressant et plus adapté en termes d'échantillonnage taxonomique. Si dans un groupe nous ne disposons que d'un seul génome fortement sujet à certains biais, nous devons quand même l'utiliser afin de représenter ce groupe dans les analyses. Par contre, si de nouveaux génomes sont disponibles, il sera possible d'une part de vérifier si ces biais sont caractéristiques de l'espèce ou de son groupe entier et de choisir un nouvel échantillonnage taxonomique plus adapté à la question scientifique posée. Par exemple, les Nanoarchaeota, groupe taxonomique d'intérêt lorsqu'on s'intéresse à la phylogénie des archées, étaient représentés jusqu'à récemment, par un seul organisme (*N. equitans*) au génome très réduit et évoluant vite. Le séquençage d'un nouveau représentant (Podar et al. 2013), ayant un génome moins réduit et évoluant moins vite, devrait permettre de réduire les biais et artefacts autrefois induits par l'utilisation de *N. equitans*.

De l'intérêt des analyses automatiques.

L'énorme masse de données générée par toutes ces séquences génomiques est bien sûr extrêmement complexe à exploiter et, par conséquent, très chronophage. Ainsi que je l'ai déjà exposé plus tôt, la plupart des analyses phylogénétiques massives n'ont, d'une part, pas réussi à dégager de nouveaux marqueurs malgré le développement de nouvelles approches pour

l'identification des marqueurs et, d'autre part, elles ont probablement été biaisées par un certain nombre d'artefacts tels que des attractions de longues branches ou la présence de paralogies cachées (Cox, Foster, Hirt, Harris, & Embley, 2008; Rinke et al., 2013...). L'analyse que j'ai menée sur l'ensemble des protéines de trois génomes complets m'a permis de mettre en évidence 200 nouveaux marqueurs, exploités immédiatement pour l'inférence de la phylogénie des archées, mais aussi un grand nombre de protéines ayant des profils évolutifs notables, avec des cas de transferts horizontaux inattendus ou des répartitions taxonomiques particulières (cf. Chapitre 3). Ces événements évolutifs n'étaient pas l'information recherchée au départ, et c'est l'observation humaine de chaque phylogénie qui a permis de les détecter et qui permettra potentiellement de les exploiter par la suite. Parallèlement, certaines protéines ont été sélectionnées comme marqueurs potentiels dans cette première étude alors qu'elles avaient des profils inattendus qui n'auraient pas été retenus par une analyse automatique dont les critères sont forcément préétablis et fixes. Par exemple, l'absence de certains ordres complets d'archées, la présence d'homologues bactériens ou eucaryotes au sein des archées issus de transferts horizontaux ponctuels et récents mais pouvant être supprimés par la suite, ou simplement l'absence de suffisamment de signal pour avoir une phylogénie individuelle bien résolue, auraient empêché la sélection automatique de ces protéines. Bien que l'ensemble des critères utilisés pour remplir les tableaux présentés en [Annexe 2](#), soit, à mon avis, parfaitement programmables un par un par un bon informaticien, il me semble beaucoup plus difficile, pour ne pas dire impossible, d'imaginer à l'avance tous les cas de figure auxquels on peut être confronté lors d'une analyse de cette ampleur. Un autre exemple provient des paralogies spécifiques à certains groupes d'archées : j'ai dû sélectionner certaines copies de paralogues, mais elles n'avaient pas forcément exactement la même répartition taxonomique ou la même longueur de branches. Les mêmes questions se sont posées sur l'analyse des jeux de données contenant les séquences bactériennes (cf. Chapitre 2). Malgré toutes ces remarques, je tiens à préciser que les critères de sélection des jeux de données, puis des séquences retenues auraient été difficilement programmables a priori et qu'ils peuvent même être considérés comme subjectifs, il n'en reste pas moins qu'ils sont reproductibles a posteriori.

Mon analyse a d'ailleurs nécessité l'automatisation de nombreuses procédures, en particulier pour réaliser les premières étapes de construction des jeux de données et des phylogénies individuelles. Cependant, la phase d'expertise des phylogénies résultantes a été cruciale. Il est important de préciser que les biais inhérents à ce protocole automatique (la limitation du nombre de séquences à 200 par jeu de données et la limite d'E-value à $e10^{-5}$ dans les recherches par BLAST, la non vérification de l'homologie des séquences dans les jeux de données, la sélection automatique des positions conservées pour la reconstruction phylogénétique et enfin, l'inférence des arbres

préliminaires en approximation de maximum de vraisemblance (méthode particulièrement sujette à l'attraction de longues branches)) ont été pris en compte dans la phase d'observation et les annotations que j'ai faites sur ces phylogénies ont été considérées avec de grandes précautions. Ceci n'aurait pas pu être le cas dans une analyse totalement automatique. Il me semble donc que l'automatisation de certaines tâches précises est nécessaire pour pouvoir traiter ces masses de données, mais qu'elle ne doit pas se faire aux dépens d'une analyse humaine à certaines étapes clés du protocole.

Le temps et les compétences humaines nécessaires à ce genre d'analyse ne sont évidemment pas compatibles avec l'analyse à flux tendu de l'ensemble de l'information génomique publiée. Un décalage se crée donc entre la disponibilité d'une information et son utilisation. De mon point de vue, il est important de choisir une question scientifique précise et de passer le temps nécessaire à l'analyse des données pour prétendre y répondre de façon fiable, plutôt que d'analyser de façon massive des données en perpétuelle expansion sans prendre en compte toutes les spécificités de la question à laquelle on veut répondre, au risque d'obtenir des résultats biaisés. La totalité des données disponibles ne sera peut-être jamais traitée, mais certaines questions scientifiques auront des réponses de plus en plus fiables.

De la multiplication des phyla archéens.

Ces dernières années (fin des années 2000 et début des années 2010), de nombreux nouveaux phyla d'archées ont été proposés. D'abord les Thaumarchaeota en 2008, définis sur l'analyse phylogénétique et génomique de *C. symbiosum* et l'analyse de diverses séquences environnementales (Brochier-Armanet et al. 2008). En 2011, Nunoura et collaborateurs proposent le phylum Aigarchaeota (Nunoura et al. 2011) à partir de l'analyse de '*Ca. Caldiarchaeum subterraneum*'. Au cours de cette dernière année, pas moins de cinq phyla ont été proposés : Geoarchaeota (Kozubal et al. 2013), suivit de Aenigmarchaeota, Diapherotrites, Nanohaloarchaeota et Parvarchaeota (Rinke et al. 2013). Les phyla Thaumarchaeota et Aigarchaeota reposent sur une analyse assez poussée des génomes disponibles et, surtout, sur le fait que ces groupes sont connus depuis longtemps par les études environnementales ; de plus, en ce qui concerne les Thaumarchaeota, des représentants cultivés étaient disponibles en 2008. Par contre, les cinq derniers nouveaux phyla reposent sur des métagénomes ou des génomes séquencés à partir de cellules uniques. Je ne mets pas en cause ici la validité de ces phyla, mais il me semble intéressant de noter le fait qu'aucune règle ne régit la taxonomie procaryote, et a fortiori des archées, en ce qui concerne les rangs taxonomiques de haut niveau (supérieurs à la classe). Certains des phyla proposés (Nanohaloarchaeota et Parvarchaeota) sont particulièrement sujets aux LBA et aucune

précaution particulière n'a été prise pour minimiser cet artefact dans l'article où Rinke et collaborateurs les ont proposés (Rinke et al. 2013). En ce qui concerne les Geoarchaeota (Kozubal et al. 2013), Aenigmarchaeota et Diapherotrites (Rinke et al. 2013), ces phyla ont été proposés dans l'article décrivant leurs génomes, mais la biologie de ces organismes reste en très grande partie inconnue. Il serait donc prudent d'attendre d'en savoir plus à leur propos et de préciser leurs positions phylogénétiques avec plus de séquences et des analyses phylogénétiques fiables avant de proposer la création de nouveaux phyla.

Du risque de découplage entre génome et organisme.

La multiplication des études de séquençage de génomes à partir de cellules uniques ou de métagénomes est extrêmement intéressante de par la quantité de nouvelles informations génomiques qu'elles apportent. Mais ces études peuvent être critiquées dans le sens où le génome n'est pas l'organisme. La séquence de tous les gènes d'un organisme ne suffit pas pour comprendre son fonctionnement cellulaire dans sa totalité. Les protéines sont prédites à partir d'ORF (« Open Reading Frame » pour « cadre ouvert de lecture ») et il n'est jamais certain que ces gènes soient réellement exprimés dans la cellule. Pour beaucoup d'ORF aucune fonction n'est assignable, et même si une fonction peut être assignée à un ORF, c'est par comparaison avec son homologue le plus proche, et dans la plupart des cas cette fonction n'est pas vérifiée de façon expérimentale. En revanche, ces génomes et métagénomes sont de nouvelles sources d'information très riches sur le métabolisme potentiel de ces organismes. Elles pourraient permettre de tenter de créer un environnement favorable à leur culture, pure ou en association avec d'autres organismes, cultures qui seraient utilisables pour des confirmations expérimentales. Dans le cas des archées, cela fut notamment le cas pour la thaumarchée *Nitrosopumilus maritimus*, la première à être isolée pour ce groupe, car des données métagénomiques avaient prédit quelques années plus tôt l'importance du métabolisme de l'ammonium chez ces organismes (Treusch et al. 2005). Il est aussi intéressant de pouvoir coupler l'analyse de génomes à d'autres techniques ne nécessitant pas d'utiliser l'organisme lui-même, telle que la transcriptomique par exemple, pour savoir si les gènes prédits sont exprimés, une idée qui est à la base des techniques de métatranscriptomiques et de métaprotéomiques de plus en plus utilisées (Shi, Tyson, and DeLong 2009; Bertin et al. 2011). L'application de ces méthodes pour comprendre le fonctionnement et l'évolution des communautés microbiennes, souvent très complexes, peut paraître utopique, mais dans la mesure où de plus en plus de génomes sont accessibles et où les organismes sont mal connus, cette voie, explorée depuis plusieurs années déjà, me semble tout à fait justifiée.

De la phylogénie moléculaire comme révolution scientifique ?

La première partie de l'Introduction de ma thèse portait principalement sur les travaux de Woese et de ses collaborateurs dans les années 1970 à 1990. Ici, je voudrais évoquer une question qui m'interpelle et l'introduire comme idée de réflexion. La notion de révolution scientifique est définie par Kuhn en 1962 dans son essai *The structure of scientific revolutions* (Kuhn 1962) comme le résultat d'un changement de paradigme dans un domaine de la science. Le paradigme scientifique selon Kuhn correspond à un modèle de pensée reposant sur une théorie, un socle d'observations, de faits et d'expériences tendant à confirmer la théorie en cours, et une série d'hypothèses et de prédictions qui viendraient confirmer la théorie. Il vient cadrer la pensée scientifique en lui donnant un socle de connaissances acquises et une méthode. Cette théorie de philosophie des sciences a été proposée d'après l'analyse de l'histoire de la physique et de la chimie principalement. En biologie, un élément fondateur est la publication de *On the Origin of Species* de Darwin en 1859 (Darwin 1859). Il a permis de donner un sens aux observations des naturalistes jusqu'à ce moment là et de donner un cadre de pensée pour l'étude des organismes vivants et des relations qui les lient. La phylogénie moléculaire, et particulièrement son application par Woese sur l'ARNr SSU à partir des années 1970, a permis de passer d'une phylogénie de l'ensemble du vivant théorique proposée par Darwin, à une phylogénie réelle du vivant. De plus, Woese a donné un nouveau cadre de pensée à la classification du vivant et à l'évolution microbienne. Le travail de Woese et de ses collaborateurs ne remet pas en question la théorie de Darwin sur l'évolution des espèces mais, au contraire, vient même la confirmer. Par contre, il apporte un nouveau cadre de pensée à un domaine de la biologie, la microbiologie, et une série d'expériences et de résultats en concordance avec ce cadre de pensée (comme la prédiction de la découverte d'une grande diversité d'archées, qui a été réalisée). Étant donné l'avancée fondamentale du travail de Woese dans les années 1970 et 1980, peut-on considérer que ce travail et cette redéfinition de la classification du vivant par la phylogénie moléculaire basée sur l'ARNr SSU soit une révolution scientifique ? Je ne prétends pas avoir de réponse à cette question et je pense qu'elle pourrait faire le sujet d'un travail à part entière.

Conclusion

Le sujet de cette thèse a été l'étude de la phylogénie des archées. Mon travail s'est articulé autour de deux axes présentés dans la partie Objectifs : la recherche de nouveaux marqueurs pour l'inférence de la phylogénie des archées et l'élucidation de la position de la racine de l'arbre des archées grâce à des homologues bactériens. Ces deux objectifs ont été atteints et les résultats obtenus font l'objet de deux articles en préparation. Grâce à l'utilisation de marqueurs et de méthodes qui n'avaient jamais été appliquées à l'étude des archées, mes travaux apportent un nouvel éclairage sur la phylogénie de ces organismes et, d'une manière plus générale, sur le travail en phylogénomique.

De façon plus précise, le Chapitre 1 a pour résultats majeurs :

- La mise en évidence de 200 nouveaux marqueurs pour la reconstruction de la phylogénie des archées, avec une grande diversité fonctionnelle par rapport aux protéines informationnelles utilisées préalablement.
- L'inférence d'une phylogénie des archées à partir de 273 protéines (200 nouvelles + 57 ribosomiques + 16 sous-unités de l'ARN polymérase et facteurs de transcription) avec un meilleur support que les phylogénies précédentes.
- La résolution de certaines relations phylogénétiques problématiques : la monophylie des méthanogènes classe I et la proposition d'une nouvelle classe les regroupant ; la non-monophylie des méthanogènes classe II et le groupement de ces méthanogènes avec les *Halobacteriales* ; la monophylie des *Thermoproteales* ; la position d'*A. saccharavorans* groupe frère d'*A. pernix* au sein des *Desulfurococcales* ; la divergence des *Thermococcales* et des *Thermoproteales* en premier au sein des euryarchées et des crenarchées respectivement ; la relation fortement soutenue entre les Thaumarchaeota et l'Aigarchaeota.

En ce qui concerne le Chapitre 2 les résultats majeurs sont :

- L'identification de 38 protéines fonctionnellement très diverses et avec des homologues bactériens utilisables pour placer la racine de l'arbre des archées, indépendamment de la présence d'homologues eucaryotes.
- L'apport d'un signal phylogénétique important par rapport aux 32 protéines ribosomiques partagées entre archées et bactéries.
- Le placement de la racine avec un support très fort entre les Euryarchaeota et les autres phyla d'archées considérés dans cette étude (Thaumarchaeota, Aigarchaeota, Crenarchaeota,

Korarchaeota).

- La proposition d'une révision de la taxonomie des archées, qui consiste principalement à déclasser quatre phyla (Thaumarchaeota, Aigarchaeota, Crenarchaeota, Korarchaeota) au rang de classe et la proposition d'un nouveau phylum englobant ces quatre groupes sous le nom de « Proteoarchaeota »

Enfin, le Chapitre 3, par l'analyse de l'histoire évolutive du « système chaperonne » DnaK-GrpE-DnaJ chez les archées, montre que des histoires évolutives complexes impliquant des événements de transfert, de duplication et de fusion de gènes ont marqué l'évolution de ce système.

Bibliographie

- Altschul, S F, W Gish, W Miller, E W Myers, and D J Lipman. 1990. "Basic Local Alignment Search Tool." *Journal of Molecular Biology* 215 (3) (October 5): 403–10. doi:10.1016/S0022-2836(05)80360-2. <http://www.ncbi.nlm.nih.gov/pubmed/2231712>.
- Altschul, S F, T L Madden, a a Schäffer, J Zhang, Z Zhang, W Miller, and D J Lipman. 1997. "Gapped BLAST and PSI-BLAST: a New Generation of Protein Database Search Programs." *Nucleic Acids Research* 25 (17) (September 1): 3389–402. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=146917&tool=pmcentrez&render type=abstract>.
- Anderson, Iain, Markus Göker, Matt Nolan, Susan Lucas, Nancy Hammon, Shweta Deshpande, Jan-Fang Cheng, et al. 2011. "Complete Genome Sequence of the Hyperthermophilic Chemolithoautotroph *Pyrolobus Fumarii* Type Strain (1A)." *Standards in Genomic Sciences* 4 (3) (July 1): 381–92. doi:10.4056/sigs.2014648. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3156397&tool=pmcentrez&render type=abstract>.
- Anderson, Iain, Reinhard Wirth, Susan Lucas, Alex Copeland, Alla Lapidus, Jan-Fang Cheng, Lynne Goodwin, et al. 2011. "Complete Genome Sequence of *Staphylothermus Hellenicus* P8." *Standards in Genomic Sciences* 5 (1) (October 15): 12–20. doi:10.4056/sigs.2054696. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3236042&tool=pmcentrez&render type=abstract>.
- Atoji, M, N Bjerrum, L Pauling, Hydrogen Bonding, A Geiger, G E Fox, E Stackebrandt, et al. 1980. "The Phylogeny of Prokaryotes." *Science (New York, N.Y.)* 209 (4455) (July 25): 457–63. <http://www.ncbi.nlm.nih.gov/pubmed/6771870>.
- Auchtung, Thomas a, Cristina D Takacs-Vesbach, and Colleen M Cavanaugh. 2006. "16S rRNA Phylogenetic Investigation of the Candidate Division 'Korarchaeota'." *Applied and Environmental Microbiology* 72 (7) (July): 5077–82. doi:10.1128/AEM.00052-06. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1489347&tool=pmcentrez&render type=abstract>.
- Baker, Brett J, Luis R Comolli, Gregory J Dick, Loren J Hauser, Doug Hyatt, Brian D Dill, Miriam L Land, Nathan C Verberkmoes, Robert L Hettich, and Jillian F Banfield. 2010. "Enigmatic, Ultrasmall, Uncultivated Archaea." *Proceedings of the National Academy of Sciences of the United States of America* 107 (19) (May 11): 8806–11. doi:10.1073/pnas.0914470107. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2889320&tool=pmcentrez&render type=abstract>.
- Baker, Brett J, Gene W Tyson, Richard I Webb, Judith Flanagan, Philip Hugenholtz, Eric E Allen, and Jillian F Banfield. 2006. "Lineages of Acidophilic Archaea Revealed by Community Genomic Analysis." *Science (New York, N.Y.)* 314 (5807) (December 22): 1933–5. doi:10.1126/science.1132690. <http://www.ncbi.nlm.nih.gov/pubmed/17185602>.
- Balch, W E, L J Magrum, G E Fox, R S Wolfe, and C R Woese. 1977. "An Ancient Divergence Among the Bacteria." *Journal of Molecular Evolution* 9 (4) (August 5): 305–11. <http://www.ncbi.nlm.nih.gov/pubmed/408502>.

- Bapteste, E, E Susko, J Leigh, D MacLeod, R L Charlebois, and W F Doolittle. 2005. "Do Orthologous Gene Phylogenies Really Support Tree-thinking?" *BMC Evolutionary Biology* 5 (January): 33. doi:10.1186/1471-2148-5-33.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1156881&tool=pmcentrez&render type=abstract>.
- Bapteste, Eric, Céline Brochier, and Yan Boucher. 2005. "Higher-level Classification of the Archaea: Evolution of Methanogenesis and Methanogens." *Archaea (Vancouver, B.C.)* 1 (5) (May): 353–63.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2685549&tool=pmcentrez&render type=abstract>.
- Barns, S M, C F Delwiche, J D Palmer, and N R Pace. 1996. "Perspectives on Archaeal Diversity, Thermophily and Monophyly from Environmental rRNA Sequences." *Proceedings of the National Academy of Sciences of the United States of America* 93 (17) (August 20): 9188–93.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=38617&tool=pmcentrez&render type=abstract>.
- Barns, S M, R E Fundyga, M W Jeffries, and N R Pace. 1994. "Remarkable Archaeal Diversity Detected in a Yellowstone National Park Hot Spring Environment." *Proceedings of the National Academy of Sciences of the United States of America* 91 (5) (March 1): 1609–13.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=43212&tool=pmcentrez&render type=abstract>.
- Benson, Dennis A, Mark Cavanaugh, Karen Clark, Ilene Karsch-Mizrachi, David J Lipman, James Ostell, and Eric W Sayers. 2013. "GenBank." *Nucleic Acids Research* 41 (Database issue) (January): D36–42. doi:10.1093/nar/gks1195.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3531190&tool=pmcentrez&render type=abstract>.
- Benson, Dennis A, Ilene Karsch-Mizrachi, Karen Clark, David J Lipman, James Ostell, and Eric W Sayers. 2012. "GenBank." *Nucleic Acids Research* 40 (Database issue) (January): D48–53. doi:10.1093/nar/gkr1202.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3245039&tool=pmcentrez&render type=abstract>.
- Bertin, Philippe N, Audrey Heinrich-Salmeron, Eric Pelletier, Florence Goulhen-Chollet, Florence Arsène-Ploetze, Sébastien Gallien, Béatrice Lauga, et al. 2011. "Metabolic Diversity Among Main Microorganisms Inside an Arsenic-rich Ecosystem Revealed by Meta- and Proteogenomics." *The ISME Journal* 5 (11) (November): 1735–47. doi:10.1038/ismej.2011.51.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3197163&tool=pmcentrez&render type=abstract>.
- Blainey, Paul C, Annika C Mosier, Anastasia Potanina, Christopher a Francis, and Stephen R Quake. 2011. "Genome of a Low-salinity Ammonia-oxidizing Archaeon Determined by Single-cell and Metagenomic Analysis." *PloS One* 6 (2) (January): e16626. doi:10.1371/journal.pone.0016626.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3043068&tool=pmcentrez&render type=abstract>.
- Boucher, Yan, W Ford Doolittle, and Peter Raynes. 2002. "Something New Under the Sea Liquid Crystal Painting" 417 (May).

- Brinkmann, H, and H Philippe. 1999. "Archaea Sister Group of Bacteria? Indications from Tree Reconstruction Artifacts in Ancient Phylogenies." *Molecular Biology and Evolution* 16 (6) (June): 817–25. <http://www.ncbi.nlm.nih.gov/pubmed/10368959>.
- Brochier, Céline, Eric Baptiste, David Moreira, and Hervé Philippe. 2002. "Eubacterial Phylogeny Based on Translational Apparatus Proteins." *Trends in Genetics : TIG* 18 (1) (January): 1–5. <http://www.ncbi.nlm.nih.gov/pubmed/11750686>.
- Brochier, Céline, Patrick Forterre, and Simonetta Gribaldo. 2004. "Archaeal Phylogeny Based on Proteins of the Transcription and Translation Machineries: Tackling the Methanopyrus Kandleri Paradox." *Genome Biology* 5 (3) (January): R17. doi:10.1186/gb-2004-5-3-r17. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=395767&tool=pmcentrez&rendertype=abstract>.
- . 2005. "An Emerging Phylogenetic Core of Archaea: Phylogenies of Transcription and Translation Machineries Converge Following Addition of New Genome Sequences." *BMC Evolutionary Biology* 5 (January): 36. doi:10.1186/1471-2148-5-36. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1177939&tool=pmcentrez&renderertype=abstract>.
- Brochier, Céline, Simonetta Gribaldo, Yvan Zivanovic, Fabrice Confalonieri, and Patrick Forterre. 2005. "Nanoarchaea: Representatives of a Novel Archaeal Phylum or a Fast-evolving Euryarchaeal Lineage Related to Thermococcales?" *Genome Biology* 6 (5) (January): R42. doi:10.1186/gb-2005-6-5-r42. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1175954&tool=pmcentrez&renderertype=abstract>.
- Brochier-Armanet, Céline, Bastien Boussau, Simonetta Gribaldo, and Patrick Forterre. 2008. "Mesophilic Crenarchaeota: Proposal for a Third Archaeal Phylum, the Thaumarchaeota." *Nature Reviews. Microbiology* 6 (3) (March): 245–52. doi:10.1038/nrmicro1852. <http://www.ncbi.nlm.nih.gov/pubmed/18274537>.
- Brochier-Armanet, Céline, Philippe Deschamps, Purificación López-García, Yvan Zivanovic, Francisco Rodríguez-Valera, and David Moreira. 2011. "Complete-fosmid and Fosmid-end Sequences Reveal Frequent Horizontal Gene Transfers in Marine Uncultured Planktonic Archaea." *The ISME Journal* 5 (8) (August): 1291–302. doi:10.1038/ismej.2011.16. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3146271&tool=pmcentrez&renderertype=abstract>.
- Brochier-Armanet, Céline, Patrick Forterre, and Simonetta Gribaldo. 2011. "Phylogeny and Evolution of the Archaea: One Hundred Genomes Later." *Current Opinion in Microbiology* 14 (3) (June): 274–81. doi:10.1016/j.mib.2011.04.015. <http://www.ncbi.nlm.nih.gov/pubmed/21632276>.
- Brock, T D, K M Brock, R T Belly, and R L Weiss. 1972. "Sulfolobus: a New Genus of Sulfur-oxidizing Bacteria Living at Low pH and High Temperature." *Archiv Für Mikrobiologie* 84 (1) (January): 54–68. <http://www.ncbi.nlm.nih.gov/pubmed/4559703>.
- Brügger, Kim, Lanming Chen, Markus Stark, Arne Zibat, Peter Redder, Andreas Ruepp, Mariana Awayez, Qunxin She, Roger a Garrett, and Hans-Peter Klenk. 2007. "The Genome of Hyperthermus Butylicus: a Sulfur-reducing, Peptide Fermenting, Neutrophilic Crenarchaeote Growing up to 108 Degrees C." *Archaea (Vancouver, B.C.)* 2 (2) (May): 127–35.

<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2686385&tool=pmcentrez&render type=abstract>.

- Bult, Carol J, Owen White, Gary J Olsen, Lixin Zhou, Robert D Fleischmann, Granger G Sutton, Judith A Blake, et al. 1996. "Complete Genome Sequence of the Methanogenic Arc Haeon , Me Thanococcus Jannaschii The Genome of Methanococcus Jannaschii" 273 (August).
- Burggraf, S, K O Stetter, P Rouviere, and C R Woese. 1991. "Methanopyrus Kandleri: An Archaeal Methanogen Unrelated to All Other Known Methanogens." *Systematic and Applied Microbiology* 14 (January): 346–51. <http://www.ncbi.nlm.nih.gov/pubmed/11540073>.
- Cavicchioli, Ricardo. 2006. "Cold-adapted Archaea." *Nature Reviews. Microbiology* 4 (5) (May): 331–43. doi:10.1038/nrmicro1390. <http://www.ncbi.nlm.nih.gov/pubmed/16715049>.
- Chen, Lanming, Kim Bru, Marie Skovgaard, Peter Redder, Qunxin She, Elfar Torarinsson, Bo Greve, et al. 2005. "The Genome of Sulfolobus Acidocaldarius , a Model Organism of the Crenarchaeota †" 187 (14): 4992–4999. doi:10.1128/JB.187.14.4992.
- Ciccarelli, Francesca D, Tobias Doerks, Christian Von Mering, Christopher J Creevey, Christian von Mering, Berend Snel, and Peer Bork. 2006. "Toward Automatic Reconstruction of a Highly Resolved Tree of Life." *Science (New York, N.Y.)* 311 (5765) (March 3): 1283–7. doi:10.1126/science.1123061. <http://www.ncbi.nlm.nih.gov/pubmed/16513982>.
- Cox, Cymon J, Peter G Foster, Robert P Hirt, Simon R Harris, and T Martin Embley. 2008. "The Archaeobacterial Origin of Eukaryotes." *Proceedings of the National Academy of Sciences of the United States of America* 105 (51) (December 23): 20356–61. doi:10.1073/pnas.0810647105. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2629343&tool=pmcentrez&render type=abstract>.
- Crisuolo, Alexis, and Simonetta Gribaldo. 2010. "BMGE (Block Mapping and Gathering with Entropy): a New Software for Selection of Phylogenetic Informative Regions from Multiple Sequence Alignments." *BMC Evolutionary Biology* 10 (January): 210. doi:10.1186/1471-2148-10-210. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3017758&tool=pmcentrez&render type=abstract>.
- Dagan, Tal, and William Martin. 2006. "The Tree of One Percent." *Genome Biology* 7 (10) (January): 118. doi:10.1186/gb-2006-7-10-118. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1794558&tool=pmcentrez&render type=abstract>.
- Darland, G, T D Brock, W Samsonoff, and S F Conti. 1970. "A Thermophilic, Acidophilic Mycoplasma Isolated from a Coal Refuse Pile." *Science (New York, N.Y.)* 170 (3965) (December 25): 1416–8. <http://www.ncbi.nlm.nih.gov/pubmed/5481857>.
- Darwin, Charles. 1859. *On the Origin of Species by Means of Natural Selection, or the Preservation of Favoured Races in the Struggle for Life*.
- De la Torre, José R, Christopher B Walker, Anitra E Ingalls, Martin Könneke, and David a Stahl. 2008. "Cultivation of a Thermophilic Ammonia Oxidizing Archaeon Synthesizing Crenarchaeol." *Environmental Microbiology* 10 (3) (March): 810–8. doi:10.1111/j.1462-2920.2007.01506.x. <http://www.ncbi.nlm.nih.gov/pubmed/18205821>.

- DeLong, E F. 1992. "Archaea in Coastal Marine Environments." *Proceedings of the National Academy of Sciences of the United States of America* 89 (12) (June 15): 5685–9. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=49357&tool=pmcentrez&rendertype=abstract>.
- DeLong, E F, K Y Wu, B B Prézelin, and R V Jovine. 1994. "High Abundance of Archaea in Antarctic Marine Picoplankton." *Nature* 371 (6499) (October 20): 695–7. doi:10.1038/371695a0. <http://www.ncbi.nlm.nih.gov/pubmed/7935813>.
- DeLong, Edward F. 2003. "Oceans of Archaea." *ASM News* 69 (10): 503–511.
- Di Giulio, Massimo. 2006. "The Non-monophyletic Origin of the tRNA Molecule and the Origin of Genes Only after the Evolutionary Stage of the Last Universal Common Ancestor (LUCA)." *Journal of Theoretical Biology* 240 (3) (June 7): 343–52. doi:10.1016/j.jtbi.2005.09.023. <http://www.ncbi.nlm.nih.gov/pubmed/16289209>.
- . 2007. "The Tree of Life Might Be Rooted in the Branch Leading to Nanoarchaeota." *Gene* 401 (1-2) (October 15): 108–13. doi:10.1016/j.gene.2007.07.004. <http://www.ncbi.nlm.nih.gov/pubmed/17689206>.
- . 2008a. "Transfer RNA Genes in Pieces Are an Ancestral Character." *EMBO Reports* 9 (9) (September): 820; author reply 820–1. doi:10.1038/embor.2008.153. http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2529356&tool=pmcentrez&render_type=abstract.
- . 2008b. "The Split Genes of Nanoarchaeum Equitans Are an Ancestral Character." *Gene* 421 (1-2) (September 15): 20–6. doi:10.1016/j.gene.2008.06.010. <http://www.ncbi.nlm.nih.gov/pubmed/18590807>.
- Do, Chuong B, Mahathi S P Mahabhashyam, Michael Brudno, and Serafim Batzoglou. 2005. "ProbCons: Probabilistic Consistency-based Multiple Sequence Alignment." *Genome Research* 15 (2) (February 1): 330–40. doi:10.1101/gr.2821705. <http://genome.cshlp.org/content/15/2/330.abstract>.
- Dobzhansky, Theodosius. 1973. "Nothing Makes Sense in Biology Except in the Light of Evolution." *The American Biology Teacher*.
- Dopson, Mark, Craig Baker-austin, Andrew Hind, John P Bowman, Philip L Bond, and A P P L E Nviron M Icrobiol. 2004. "Characterization of Ferroplasma Isolates and Ferroplasma Acidarmanus Sp. Nov., Extreme Acidophiles from Acid Mine Drainage and Industrial Bioleaching Environments" 70 (4): 2079–2088. doi:10.1128/AEM.70.4.2079.
- Dorn, Karolin V, Felix Willmund, Christian Schwarz, Christine Henselmann, Thomas Pohl, Barbara Hess, Daniel Veyel, et al. 2010. "Chloroplast DnaJ-like Proteins 3 and 4 (CDJ3/4) from Chlamydomonas Reinhardtii Contain Redox-active Fe-S Clusters and Interact with Stromal HSP70B." *The Biochemical Journal* 427 (2) (April 15): 205–15. doi:10.1042/BJ20091412. <http://www.ncbi.nlm.nih.gov/pubmed/20113313>.
- Edgar, Robert C. 2004. "MUSCLE: Multiple Sequence Alignment with High Accuracy and High Throughput." *Nucleic Acids Research* 32 (5) (January): 1792–7. doi:10.1093/nar/gkh340. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=390337&tool=pmcentrez&rendertype=abstract>.

- Elkins, James G, Mircea Podar, David E Graham, Kira S Makarova, Yuri Wolf, Lennart Randau, Brian P Hedlund, et al. 2008. "A Korarchaeal Genome Reveals Insights into the Evolution of the Archaea." *Proceedings of the National Academy of Sciences of the United States of America* 105 (23) (June 10): 8102–7. doi:10.1073/pnas.0801980105. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2430366&tool=pmcentrez&render type=abstract>.
- Eme, Laura, Laila J Reigstad, Anja Spang, Anders Lanzén, Thomas Weinmaier, Thomas Rattei, Christa Schleper, and Céline Brochier-Armanet. 2013. "Metagenomics of Kamchatkan Hot Spring Filaments Reveal Two New Major (hyper)thermophilic Lineages Related to Thaumarchaeota." *Research in Microbiology* 164 (5) (June): 425–38. doi:10.1016/j.resmic.2013.02.006. <http://www.ncbi.nlm.nih.gov/pubmed/23470515>.
- Farlow, W. G. 1880. "On the Nature of the Peculiar Reddening of Salted Codfish During the Summer Season." *Fish Comm.* 3: 969–97.
- Fitch, W M, and E Margoliash. 1967. "Construction of Phylogenetic Trees." *Science (New York, N.Y.)* 155 (3760) (January 20): 279–84. <http://www.ncbi.nlm.nih.gov/pubmed/5334057>.
- Foster, Peter G, Cymon J Cox, and T Martin Embley. 2009. "The Primary Divisions of Life: a Phylogenomic Approach Employing Composition-heterogeneous Methods." *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences* 364 (1527) (August 12): 2197–207. doi:10.1098/rstb.2009.0034. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2873002&tool=pmcentrez&render type=abstract>.
- Fox, G. E., K. R. Pechman, and C. R. Woese. 1977. "Comparative Cataloging of 16S Ribosomal Ribonucleic Acid: Molecular Approach to Procaryotic Systematics." *International Journal of Systematic Bacteriology* 27 (1) (January 1): 44–57. doi:10.1099/00207713-27-1-44. <http://ijs.sgmjournals.org/cgi/doi/10.1099/00207713-27-1-44>.
- Fox, George E, Linda J Magrum, W E Balch, R S Wolfe, Carl R Woese, William E Balcht, and Ralph S Wolfef. 1977. "Classification of Methanogenic Bacteria by 16S Ribosomal RNA Characterization." *Proceedings of the National Academy of Sciences of the United States of America* 74 (10) (October): 4537–41. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=431980&tool=pmcentrez&rendert ype=abstract>.
- Franzmann, P D, Y Liu, D L Balkwill, H C Aldrich, E Conway de Macario, and D R Boone. 1997. "Methanogenium Frigidum Sp. Nov., a Psychrophilic, H₂-using Methanogen from Ace Lake, Antarctica." *International Journal of Systematic Bacteriology* 47 (4) (October): 1068–72. <http://www.ncbi.nlm.nih.gov/pubmed/9336907>.
- Franzmann, P.D., N. Springer, W. Ludwig, E. Conway De Macario, and M. Rohde. 1992. "A Methanogenic Archaeon from Ace Lake, Antarctica: Methanococcoides Burtonii Sp. Nov." *Systematic and Applied Microbiology* 15 (4) (December): 573–581. doi:10.1016/S0723-2020(11)80117-7. [http://dx.doi.org/10.1016/S0723-2020\(11\)80117-7](http://dx.doi.org/10.1016/S0723-2020(11)80117-7).
- Franzmann, P.D., E. Stackebrandt, K. Sanderson, J.K. Volkman, D.E. Cameron, P.L. Stevenson, T.A. Mcmeekin, and H.R. Burton. 1988. "Halobacterium Lacusprofundi Sp. Nov., a Halophilic Bacterium Isolated from Deep Lake, Antarctica." *Systematic and Applied Microbiology* 11 (1) (November): 20–27. doi:10.1016/S0723-2020(88)80044-4. [http://dx.doi.org/10.1016/S0723-2020\(88\)80044-4](http://dx.doi.org/10.1016/S0723-2020(88)80044-4).

- French, Sarah L, Thomas J Santangelo, Ann L Beyer, and John N Reeve. 2007. "Transcription and Translation Are Coupled in Archaea." *Molecular Biology and Evolution* 24 (4) (April): 893–5. doi:10.1093/molbev/msm007. <http://www.ncbi.nlm.nih.gov/pubmed/17237472>.
- Fuhrman, J a, K McCallum, and a a Davis. 1992. "Novel Major Archaeobacterial Group from Marine Plankton." *Nature* 356 (6365) (March 12): 148–9. doi:10.1038/356148a0. <http://www.ncbi.nlm.nih.gov/pubmed/1545865>.
- Ghai, Rohit, Lejla Pašić, Ana Beatriz Fernández, Ana-Belen Martin-Cuadrado, Carolina Megumi Mizuno, Katherine D McMahon, R Thane Papke, et al. 2011. "New Abundant Microbial Groups in Aquatic Hypersaline Environments." *Scientific Reports* 1 (January): 135. doi:10.1038/srep00135. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3216616&tool=pmcentrez&render type=abstract>.
- Gogarten, J P, H Kibak, P Dittrich, L Taiz, E J Bowman, B J Bowman, M F Manolson, R J Poole, T Date, and T Oshima. 1989. "Evolution of the Vacuolar H⁺-ATPase: Implications for the Origin of Eukaryotes." *Proceedings of the National Academy of Sciences of the United States of America* 86 (17) (September): 6661–5. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=297905&tool=pmcentrez&rendert ype=abstract>.
- Göker, Markus, Brittany Held, Alla Lapidus, Matt Nolan, Stefan Spring, Montri Yasawong, Susan Lucas, et al. 2010. "Complete Genome Sequence of Ignisphaera Aggregans Type Strain (AQ1.S1)." *Standards in Genomic Sciences* 3 (1) (January): 66–75. doi:10.4056/sigs.1072907. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3035270&tool=pmcentrez&render type=abstract>.
- Gribaldo, S, and P Cammarano. 1998. "The Root of the Universal Tree of Life Inferred from Anciently Duplicated Genes Encoding Components of the Protein-targeting Machinery." *Journal of Molecular Evolution* 47 (5) (November): 508–16. <http://www.ncbi.nlm.nih.gov/pubmed/9797401>.
- Gribaldo, Simonetta, and Celine Brochier. 2009. "Phylogeny of Prokaryotes: Does It Exist and Why Should We Care?" *Research in Microbiology* 160 (7) (September): 513–21. doi:10.1016/j.resmic.2009.07.006. <http://www.ncbi.nlm.nih.gov/pubmed/19631737>.
- Gribaldo, Simonetta, and Celine Brochier-Armanet. 2006. "The Origin and Evolution of Archaea: a State of the Art." *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences* 361 (1470) (June 29): 1007–22. doi:10.1098/rstb.2006.1841. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1578729&tool=pmcentrez&render type=abstract>.
- Gribaldo, Simonetta, and Céline Brochier-Armanet. 2012. "Time for Order in Microbial Systematics." *Trends in Microbiology* 20 (5) (May): 209–10. doi:10.1016/j.tim.2012.02.006. <http://www.ncbi.nlm.nih.gov/pubmed/22440793>.
- Gribaldo, Simonetta, Anthony M Poole, Vincent Daubin, Patrick Forterre, and Céline Brochier-Armanet. 2010. "The Origin of Eukaryotes and Their Relationship with the Archaea: Are We at a Phylogenomic Impasse?" *Nature Reviews. Microbiology* 8 (10) (October): 743–52. doi:10.1038/nrmicro2426. <http://www.ncbi.nlm.nih.gov/pubmed/20844558>.
- Guindon, Stéphane, Jean-François Dufayard, Vincent Lefort, Maria Anisimova, Wim Hordijk, and Olivier Gascuel. 2010. "New Algorithms and Methods to Estimate Maximum-likelihood

Phylogenies: Assessing the Performance of PhyML 3.0.” *Systematic Biology* 59 (3) (May): 307–21. doi:10.1093/sysbio/syq010. <http://www.ncbi.nlm.nih.gov/pubmed/20525638>.

Gupta, R, J M Lanter, and C R Woese. 1983. “Sequence of the 16S Ribosomal RNA from *Halobacterium Volcanii*, an Archaeobacterium.” *Science (New York, N.Y.)* 221 (4611) (August 12): 656–9. doi:10.1126/science.221.4611.656. <http://www.ncbi.nlm.nih.gov/pubmed/17787735>.

Gupta, R S. 2000. “The Natural Evolutionary Relationships Among Prokaryotes.” *Critical Reviews in Microbiology* 26 (2) (January): 111–31. doi:10.1080/10408410091154219. <http://www.ncbi.nlm.nih.gov/pubmed/10890353>.

Guy, Lionel, and Thijs J G Ettema. 2011. “The Archaeal ‘TACK’ Superphylum and the Origin of Eukaryotes.” *Trends in Microbiology* 19 (12) (December): 580–7. doi:10.1016/j.tim.2011.09.002. <http://www.ncbi.nlm.nih.gov/pubmed/22018741>.

Hafenbradl, Doris, Martin Keller, Reinhard Dirmeier, Reinhard Rachel, Petra Roßnagel, Siegfried Burggraf, Harald Huber, and Karl O Stetter. 1996. “a Novel Hyperthermophilic Archaeum That Oxidizes Fe 2 + at Neutral pH Under Anoxic Conditions” 2: 308–314.

Hallam, Steven J, Konstantinos T Konstantinidis, Nik Putnam, Christa Schleper, Yoh-ichi Watanabe, Junichi Sugahara, Christina Preston, José de la Torre, Paul M Richardson, and Edward F DeLong. 2006. “Genomic Analysis of the Uncultivated Marine Crenarchaeote *Cenarchaeum Symbiosum*.” *Proceedings of the National Academy of Sciences of the United States of America* 103 (48) (November 28): 18296–301. doi:10.1073/pnas.0608549103. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1643844&tool=pmcentrez&render type=abstract>.

Harris, J Kirk, Scott T Kelley, George B Spiegelman, and Norman R Pace. 2003. “The Genetic Core of the Universal Ancestor” (February): 407–412. doi:10.1101/gr.652803.

Hartmann, R, H D Sickinger, and D Oesterhelt. 1980. “Anaerobic Growth of Halobacteria.” *Proceedings of the National Academy of Sciences of the United States of America* 77 (7) (July): 3821–5. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=349718&tool=pmcentrez&render type=abstract>.

Hatzenpichler, Roland, Elena V Lebedeva, Eva Spieck, Kilian Stoecker, Andreas Richter, Holger Daims, and Michael Wagner. 2008. “A Moderately Thermophilic Ammonia-oxidizing Crenarchaeote from a Hot Spring.” *Proceedings of the National Academy of Sciences of the United States of America* 105 (6) (February 12): 2134–9. doi:10.1073/pnas.0708857105. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2538889&tool=pmcentrez&render type=abstract>.

Huber, Harald, Michael J Hohn, Reinhard Rachel, Tanja Fuchs, Verena C Wimmer, and Karl O Stetter. 2002. “A New Phylum of Archaea Represented by a Nanosized Hyperthermophilic Symbiont.” *Nature* 417 (6884) (May 2): 63–7. doi:10.1038/417063a. <http://www.ncbi.nlm.nih.gov/pubmed/11986665>.

Huber, R, M Sacher, a Vollmann, H Huber, and D Rose. 2000. “Respiration of Arsenate and Selenate by Hyperthermophilic Archaea.” *Systematic and Applied Microbiology* 23 (3) (October): 305–14. doi:10.1016/S0723-2020(00)80058-2. <http://www.ncbi.nlm.nih.gov/pubmed/11108007>.

- Iverson, Vaughn, Robert M Morris, Christian D Frazar, Chris T Berthiaume, Rhonda L Morales, and E Virginia Armbrust. 2012. "Untangling Genomes from Metagenomes: Revealing an Uncultured Class of Marine Euryarchaeota." *Science (New York, N.Y.)* 335 (6068) (February 3): 587–90. doi:10.1126/science.1212665. <http://www.ncbi.nlm.nih.gov/pubmed/22301318>.
- Iwabe, N, K Kuma, M Hasegawa, S Osawa, and T Miyata. 1989. "Evolutionary Relationship of Archaeobacteria, Eubacteria, and Eukaryotes Inferred from Phylogenetic Trees of Duplicated Genes." *Proceedings of the National Academy of Sciences of the United States of America* 86 (23) (December): 9355–9. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=298494&tool=pmcentrez&rendertype=abstract>.
- Jobb, Gangolf, Arndt von Haeseler, and Korbinian Strimmer. 2004. "TREEFINDER: a Powerful Graphical Analysis Environment for Molecular Phylogenetics." *BMC Evolutionary Biology* 4 (June 28): 18. doi:10.1186/1471-2148-4-18. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=459214&tool=pmcentrez&rendertype=abstract>.
- Joulian, C, B K Patel, B Ollivier, J L Garcia, and P a Roger. 2000. "Methanobacterium Oryzae Sp. Nov., a Novel Methanogenic Rod Isolated from a Philippines Ricefield." *International Journal of Systematic and Evolutionary Microbiology* 50 Pt 2 (March): 525–8. <http://www.ncbi.nlm.nih.gov/pubmed/10758856>.
- Karner, M B, E F DeLong, and D M Karl. 2001. "Archaeal Dominance in the Mesopelagic Zone of the Pacific Ocean." *Nature* 409 (6819) (January 25): 507–10. doi:10.1038/35054051. <http://www.ncbi.nlm.nih.gov/pubmed/11206545>.
- Katoh, Kazutaka, Kei-ichi Kuma, Hiroyuki Toh, and Takashi Miyata. 2005. "MAFFT Version 5: Improvement in Accuracy of Multiple Sequence Alignment." *Nucleic Acids Research* 33 (2) (January): 511–8. doi:10.1093/nar/gki198. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=548345&tool=pmcentrez&rendertype=abstract>.
- Klenk, H P, and W Zillig. 1994. "DNA-dependent RNA Polymerase Subunit B as a Tool for Phylogenetic Reconstructions: Branching Topology of the Archaeal Domain." *Journal of Molecular Evolution* 38 (4) (April): 420–32. <http://www.ncbi.nlm.nih.gov/pubmed/8007009>.
- Knittel, Katrin, and Antje Boetius. 2009. "Anaerobic Oxidation of Methane: Progress with an Unknown Process." *Annual Review of Microbiology* 63 (January): 311–34. doi:10.1146/annurev.micro.61.080706.093130. <http://www.ncbi.nlm.nih.gov/pubmed/19575572>.
- KOCUR, M., and W. Hodgkiss. 1973. "Taxonomic Status of the Genus Halococcus Schoop." *International Journal of Systematic Bacteriology* 23 (2) (April 1): 151–156. doi:10.1099/00207713-23-2-151. <http://ijs.sgmjournals.org/content/23/2/151.abstract>.
- Könneke, Martin, Anne E Bernhard, José R de la Torre, Christopher B Walker, John B Waterbury, and David a Stahl. 2005. "Isolation of an Autotrophic Ammonia-oxidizing Marine Archaeon." *Nature* 437 (7058) (September 22): 543–6. doi:10.1038/nature03911. <http://www.ncbi.nlm.nih.gov/pubmed/16177789>.
- Kostka, Martin, Magdalena Uzlikova, Ivan Cepicka, and Jaroslav Flegr. 2008. "SlowFaster, a User-friendly Program for Slow-fast Analysis and Its Application on Phylogeny of Blastocystis." *BMC Bioinformatics* 9 (January): 341. doi:10.1186/1471-2105-9-341.

<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2529323&tool=pmcentrez&render type=abstract>.

- Kozubal, Mark a, Margaret Romine, Ryan deM Jennings, Zack J Jay, Susannah G Tringe, Doug B Rusch, Jacob P Beam, Lee Ann McCue, and William P Inskeep. 2013. "Geoarchaeota: a New Candidate Phylum in the Archaea from High-temperature Acidic Iron Mats in Yellowstone National Park." *The ISME Journal* 7 (3) (March): 622–34. doi:10.1038/ismej.2012.132. <http://www.ncbi.nlm.nih.gov/pubmed/23151644>.
- Kuhn, Thomas S. 1962. *The Structure of Scientific Revolutions*.
- Kurr, Margit, Robert Huber, Helmut König, Holger W. Jannasch, Hans Fricke, Antonio Trincone, Jakob K. Kristjansson, and Karl O. Stetter. 1991. "Methanopyrus Kandleri, Gen. and Sp. Nov. Represents a Novel Group of Hyperthermophilic Methanogens, Growing at 110°C." *Archives of Microbiology* 156 (4) (September): 239–247. doi:10.1007/BF00262992. <http://link.springer.com/10.1007/BF00262992>.
- Kyrpides, N C. 1999. "Genomes OnLine Database (GOLD 1.0): a Monitor of Complete and Ongoing Genome Projects World-wide." *Bioinformatics (Oxford, England)* 15 (9) (September): 773–4. <http://www.ncbi.nlm.nih.gov/pubmed/10498782>.
- Larkin, M A, G Blackshields, N P Brown, R Chenna, P A McGettigan, H McWilliam, F Valentin, et al. 2007. "Clustal W and Clustal X Version 2.0." *Bioinformatics (Oxford, England)* 23 (21) (November 1): 2947–8. doi:10.1093/bioinformatics/btm404. <http://www.ncbi.nlm.nih.gov/pubmed/17846036>.
- Lartillot, Nicolas, Thomas Lepage, and Samuel Blanquart. 2009. "PhyloBayes 3: a Bayesian Software Package for Phylogenetic Reconstruction and Molecular Dating." *Bioinformatics (Oxford, England)* 25 (17) (September 1): 2286–8. doi:10.1093/bioinformatics/btp368. <http://www.ncbi.nlm.nih.gov/pubmed/19535536>.
- Lasek-Nesselquist, Erica, and Johann Peter Gogarten. 2013. "The Effects of Model Choice and Mitigating Bias on the Ribosomal Tree of Life." *Molecular Phylogenetics and Evolution* 69 (1) (May 21): 17–38. doi:10.1016/j.ympev.2013.05.006. <http://www.ncbi.nlm.nih.gov/pubmed/23707703>.
- Le, Si Quang, and Olivier Gascuel. 2008. "An Improved General Amino Acid Replacement Matrix." *Molecular Biology and Evolution* 25 (7) (July): 1307–20. doi:10.1093/molbev/msn067. <http://www.ncbi.nlm.nih.gov/pubmed/18367465>.
- Leigh, John A, Sonja-Verena Albers, Haruyuki Atomi, and Thorsten Allers. 2011. "Model Organisms for Genetics in the Domain Archaea: Methanogens, Halophiles, Thermococcales and Sulfolobales." *FEMS Microbiology Reviews* 35 (4) (July): 577–608. doi:10.1111/j.1574-6976.2011.00265.x. <http://www.ncbi.nlm.nih.gov/pubmed/21265868>.
- Liu, Li-Jun, Xiao-Yan You, Xu Guo, Shuang-Jiang Liu, and Cheng-Ying Jiang. 2011. "Metallosphaera Cuprina Sp. Nov., an Acidothermophilic, Metal-mobilizing Archaeon." *International Journal of Systematic and Evolutionary Microbiology* 61 (Pt 10) (October): 2395–400. doi:10.1099/ijs.0.026591-0. <http://www.ncbi.nlm.nih.gov/pubmed/21057050>.
- Liu, Yuchen, and William B Whitman. 2008. "Metabolic, Phylogenetic, and Ecological Diversity of the Methanogenic Archaea." *Annals of the New York Academy of Sciences* 1125 (March): 171–89. doi:10.1196/annals.1419.019. <http://www.ncbi.nlm.nih.gov/pubmed/18378594>.

- López-garcía, Purificación, Céline Brochier, David Moreira, and Francisco Rodríguez-valera. 2004. “Comparative Analysis of a Genome Fragment of an Uncultivated Mesopelagic Crenarchaeote Reveals Multiple Horizontal Gene Transfers” 6: 19–34. doi:10.1046/j.1462-2920.2003.00533.x.
- Magrum, L J, K R Luehrsen, and C R Woese. 1978. “Are Extreme Halophiles Actually ‘Bacteria’?” *Journal of Molecular Evolution* 11 (1) (May 12): 1–8. <http://www.ncbi.nlm.nih.gov/pubmed/660662>.
- Mardanov, Andrey V, Nikolai V Ravin, Vitali a Svetlitchnyi, Alexey V Beletsky, Margarita L Miroshnichenko, Elizaveta a Bonch-Osmolovskaya, and Konstantin G Skryabin. 2009. “Metabolic Versatility and Indigenous Origin of the Archaeon *Thermococcus Sibiricus*, Isolated from a Siberian Oil Reservoir, as Revealed by Genome Analysis.” *Applied and Environmental Microbiology* 75 (13) (July): 4580–8. doi:10.1128/AEM.00718-09. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2704819&tool=pmcentrez&render type=abstract>.
- Margulis, L, and R Guerrero. 1991. “Kingdoms in Turmoil.” *New Scientist (1971)* 1761 (March 23): 46–50. <http://www.ncbi.nlm.nih.gov/pubmed/11538109>.
- Matte-Tailliez, Oriane, Céline Brochier, Patrick Forterre, and Hervé Philippe. 2002. “Archaeal Phylogeny Based on Ribosomal Proteins.” *Molecular Biology and Evolution* 19 (5) (May): 631–9. <http://www.ncbi.nlm.nih.gov/pubmed/11961097>.
- Mayr, Ernst. 1990. “A Natural System of Organisms.” *Nature* 348 (6301) (December 6): 491–491. doi:10.1038/348491a0. <http://dx.doi.org/10.1038/348491a0>.
- . 1991. “More Natural Classification.” *Nature* 353 (6340) (September 12): 122–122. doi:10.1038/353122a0. <http://www.nature.com/nature/journal/v353/n6340/pdf/353122a0.pdf>.
- Mihajlovski, Agnès, Monique Alric, and Jean-François Brugère. 2008. “A Putative New Order of Methanogenic Archaea Inhabiting the Human Gut, as Revealed by Molecular Analyses of the *mcrA* Gene.” *Research in Microbiology* 159 (7-8): 516–21. doi:10.1016/j.resmic.2008.06.007. <http://www.ncbi.nlm.nih.gov/pubmed/18644435>.
- Mills, Heath J, Robert J Martinez, Sandra Story, and Patricia A Sobecky. 2005. “Characterization of Microbial Community Structure in Gulf of Mexico Gas Hydrates : Comparative Analysis of DNA- and RNA-Derived Clone Libraries” 71 (6): 3235–3247. doi:10.1128/AEM.71.6.3235.
- Moore, R L, and B J McCarthy. 1967. “Comparative Study of Ribosomal Ribonucleic Acid Cistrons in Enterobacteria and Myxobacteria.” *Journal of Bacteriology* 94 (4) (October): 1066–74. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=276777&tool=pmcentrez&render type=abstract>.
- Mosier, Annika C, Eric E Allen, Maria Kim, Steven Ferriera, and Christopher a Francis. 2012. “Genome Sequence of ‘Candidatus Nitrosopumilus Salaria’ BD31, an Ammonia-oxidizing Archaeon from the San Francisco Bay Estuary.” *Journal of Bacteriology* 194 (8) (April): 2121–2. doi:10.1128/JB.00013-12. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3318490&tool=pmcentrez&render type=abstract>.

- Mullakhanbhai, M F, and H Larsen. 1975. "Halobacterium Volcanii Spec. Nov., a Dead Sea Halobacterium with a Moderate Salt Requirement." *Archives of Microbiology* 104 (3) (August 28): 207–14. <http://www.ncbi.nlm.nih.gov/pubmed/1190944>.
- Narasingarao, Priya, Sheila Podell, Juan a Ugalde, Céline Brochier-Armanet, Joanne B Emerson, Jochen J Brocks, Karla B Heidelberg, Jillian F Banfield, and Eric E Allen. 2012. "De Novo Metagenomic Assembly Reveals Abundant Novel Major Lineage of Archaea in Hypersaline Microbial Communities." *The ISME Journal* 6 (1) (January): 81–93. doi:10.1038/ismej.2011.78. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3246234&tool=pmcentrez&render type=abstract>.
- Notredame, C, D G Higgins, and J Heringa. 2000. "T-Coffee: A Novel Method for Fast and Accurate Multiple Sequence Alignment." *Journal of Molecular Biology* 302 (1) (September 8): 205–17. doi:10.1006/jmbi.2000.4042. <http://www.ncbi.nlm.nih.gov/pubmed/10964570>.
- Nunoura, Takuro, Hisako Hirayama, Hideto Takami, Hanako Oida, Shinro Nishi, Shigeru Shimamura, Yohey Suzuki, et al. 2005. "Genetic and Functional Properties of Uncultivated Thermophilic Crenarchaeotes from a Subsurface Gold Mine as Revealed by Analysis of Genome Fragments." *Environmental Microbiology* 7 (12) (December): 1967–84. doi:10.1111/j.1462-2920.2005.00881.x. <http://www.ncbi.nlm.nih.gov/pubmed/16309394>.
- Nunoura, Takuro, Hanako Oida, Miwako Nakaseama, Ayako Kosaka, Satoru B Ohkubo, Toru Kikuchi, Hiromi Kazama, et al. 2010. "Archaeal Diversity and Distribution Along Thermal and Geochemical Gradients in Hydrothermal Sediments at the Yonaguni Knoll IV Hydrothermal Field in the Southern Okinawa Trough." *Applied and Environmental Microbiology* 76 (4) (February): 1198–211. doi:10.1128/AEM.00924-09. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2820965&tool=pmcentrez&render type=abstract>.
- Nunoura, Takuro, Yoshihiro Takaki, Jungo Kakuta, Shinro Nishi, Junichi Sugahara, Hiromi Kazama, Gab-Joo Chee, et al. 2011. "Insights into the Evolution of Archaea and Eukaryotic Protein Modifier Systems Revealed by the Genome of a Novel Archaeal Group." *Nucleic Acids Research* 39 (8) (April): 3204–23. doi:10.1093/nar/gkq1228. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3082918&tool=pmcentrez&render type=abstract>.
- Offre, Pierre, Anja Spang, and Christa Schleper. 2013. "Archaea in Biogeochemical Cycles." *Annual Review of Microbiology* 67 (1) (November): 130628184403000. doi:10.1146/annurev-micro-092412-155614. <http://www.annualreviews.org/doi/abs/10.1146/annurev-micro-092412-155614>.
- Olsen, G J, D J Lane, S J Giovannoni, N R Pace, and D a Stahl. 1986. "Microbial Ecology and Evolution: a Ribosomal RNA Approach." *Annual Review of Microbiology* 40 (January): 337–65. doi:10.1146/annurev.mi.40.100186.002005. <http://www.ncbi.nlm.nih.gov/pubmed/2430518>.
- Olsen, G J, N R Pace, M Nuell, B P Kaine, R Gupta, and C R Woese. 1985. "Sequence of the 16S rRNA Gene from the Thermoacidophilic Archaeobacterium Sulfolobus Solfataricus and Its Evolutionary Implications." *Journal of Molecular Evolution* 22 (4) (January): 301–7. <http://www.ncbi.nlm.nih.gov/pubmed/3936935>.

- Olsen, G J, and C R Woese. 1993. "Ribosomal RNA: a Key to Phylogeny." *FASEB Journal : Official Publication of the Federation of American Societies for Experimental Biology* 7 (1) (January): 113–23. <http://www.ncbi.nlm.nih.gov/pubmed/8422957>.
- Pace, B, and L L Campbell. 1971. "Homology of Ribosomal Ribonucleic Acid Diverse Bacterial Species with Escherichia Coli and Bacillus Stearothermophilus." *Journal of Bacteriology* 107 (2) (August): 543–7. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=246958&tool=pmcentrez&rendertype=abstract>.
- Pace, N R, G J Olsen, and C R Woese. 1986. "Ribosomal RNA Phylogeny and the Primary Lines of Evolutionary Descent." *Cell* 45 (3) (May 9): 325–6. <http://www.ncbi.nlm.nih.gov/pubmed/3084106>.
- Park, Soo-Je, Jong-Geol Kim, Man-Young Jung, So-Jeong Kim, In-Tae Cha, Rohit Ghai, Ana-Belén Martín-Cuadrado, Francisco Rodríguez-Valera, and Sung-Keun Rhee. 2012. "Draft Genome Sequence of an Ammonia-oxidizing Archaeon, 'Candidatus Nitrosopumilus Sediminis' AR2, from Svalbard in the Arctic Circle." *Journal of Bacteriology* 194 (24) (December): 6948–9. doi:10.1128/JB.01869-12. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3510607&tool=pmcentrez&rendertype=abstract>.
- Paul, Kristina, James O Nonoh, Lena Mikulski, and Andreas Brune. 2012. "'Methanoplasmatales,' Thermoplasmatales-related Archaea in Termite Guts and Other Environments, Are the Seventh Order of Methanogens." *Applied and Environmental Microbiology* 78 (23) (December): 8245–53. doi:10.1128/AEM.02193-12. <http://www.ncbi.nlm.nih.gov/pubmed/23001661>.
- Perevalova, Anna a, Salima Kh Bidzhieva, Ilya V Kublanov, Kai-Uwe Hinrichs, Xiaolei L Liu, Andrey V Mardanov, Alexander V Lebedinsky, and Elizaveta a Bonch-Osmolovskaya. 2010. "Fervidicoccus Fontis Gen. Nov., Sp. Nov., an Anaerobic, Thermophilic Crenarchaeote from Terrestrial Hot Springs, and Proposal of Fervidicocaceae Fam. Nov. and Fervidicoccales Ord. Nov." *International Journal of Systematic and Evolutionary Microbiology* 60 (Pt 9) (September): 2082–8. doi:10.1099/ijs.0.019042-0. <http://www.ncbi.nlm.nih.gov/pubmed/19837732>.
- Petitjean, Céline, David Moreira, Purificación López-García, and Céline Brochier-Armanet. 2012. "Horizontal Gene Transfer of a Chloroplast DnaJ-Fer Protein to Thaumarchaeota and the Evolutionary History of the DnaK Chaperone System in Archaea." *BMC Evolutionary Biology* 12 (January): 226. doi:10.1186/1471-2148-12-226. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3564930&tool=pmcentrez&rendertype=abstract>.
- Philippe, H. 1993. "MUST, a Computer Package of Management Utilities for Sequences and Trees." *Nucleic Acids Research* 21 (22) (November 11): 5264–72. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=310646&tool=pmcentrez&rendertype=abstract>.
- Philippe, H, and J Laurent. 1998. "How Good Are Deep Phylogenetic Trees?" *Current Opinion in Genetics & Development* 8 (6) (December): 616–23. <http://www.ncbi.nlm.nih.gov/pubmed/9914208>.

- Pisani, Davide, James a Cotton, and James O McInerney. 2007. "Supertrees Disentangle the Chimerical Origin of Eukaryotic Genomes." *Molecular Biology and Evolution* 24 (8) (August): 1752–60. doi:10.1093/molbev/msm095. <http://www.ncbi.nlm.nih.gov/pubmed/17504772>.
- Podar, Mircea, Iain Anderson, Kira S Makarova, James G Elkins, Natalia Ivanova, Mark a Wall, Athanasios Lykidis, et al. 2008. "A Genomic Analysis of the Archaeal System *Ignicoccus hospitalis*-*Nanoarchaeum Equitans*." *Genome Biology* 9 (11) (January): R158. doi:10.1186/gb-2008-9-11-r158. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2614490&tool=pmcentrez&render type=abstract>.
- Podar, Mircea, Kira S Makarova, David E Graham, Yuri I Wolf, Eugene V Koonin, and Anna-Louise Reysenbach. 2013. "Insights into Archaeal Evolution and Symbiosis from the Genomes of a Nanoarchaeon and Its Inferred Crenarchaeal Host from Obsidian Pool, Yellowstone National Park." *Biology Direct* 8 (January): 9. doi:10.1186/1745-6150-8-9. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3655853&tool=pmcentrez&render type=abstract>.
- Preston, Christina M, K E Ying Wu, Tadeusz F Molinskit, and Edward F Delong. 1996. "A Psychrophilic Crenarchaeon Inhabits a Marine Sponge :." 93 (June): 6241–6246.
- Price, Morgan N, Paramvir S Dehal, and Adam P Arkin. 2010. "FastTree 2--approximately Maximum-likelihood Trees for Large Alignments." Edited by Art F. Y. Poon. *PloS One* 5 (3) (January): e9490. doi:10.1371/journal.pone.0009490. <http://dx.plos.org/10.1371/journal.pone.0009490>.
- Prokofeva, M I, N a Kostrikina, T V Kolganova, T P Tourova, a M Lysenko, a V Lebedinsky, and E a Bonch-Osmolovskaya. 2009. "Isolation of the Anaerobic Thermoacidophilic Crenarchaeote *Acidilobus Saccharovorans* Sp. Nov. and Proposal of Acidilobales Ord. Nov., Including Acidilobaceae Fam. Nov. and Caldisphaeraceae Fam. Nov." *International Journal of Systematic and Evolutionary Microbiology* 59 (Pt 12) (December): 3116–22. doi:10.1099/ijs.0.010355-0. <http://www.ncbi.nlm.nih.gov/pubmed/19643887>.
- Rachel, Reinhard, Irith Wyszchony, Sabine Riehl, and Harald Huber. 2002. "The Ultrastructure of *Ignicoccus*: Evidence for a Novel Outer Membrane and for Intracellular Vesicle Budding in an Archaeon." *Archaea (Vancouver, B.C.)* 1 (1) (March): 9–18. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2685547&tool=pmcentrez&render type=abstract>.
- Reeburgh, William S. 2007. "Oceanic Methane Biogeochemistry." *Chemical Reviews* 107 (2) (February): 486–513. doi:10.1021/cr050362v. <http://www.ncbi.nlm.nih.gov/pubmed/17261072>.
- Reno, Michael L, Nicole L Held, Christopher J Fields, Patricia V Burke, and Rachel J Whitaker. 2009. "Biogeography of the *Sulfolobus Islandicus* Pan-genome." *Proceedings of the National Academy of Sciences of the United States of America* 106 (21) (May 26): 8605–10. doi:10.1073/pnas.0808945106. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2689034&tool=pmcentrez&render type=abstract>.
- Reykjaví, Is-, V T Marteinson, S Hauksdóttir, C F Hobel, H Kristmannsdóttir, G O Hreggvidsson, and J K Kristjánsson. 2001. "Phylogenetic Diversity Analysis of Subterranean Hot Springs in Iceland." *Applied and Environmental Microbiology* 67 (9) (September): 4242–8.

doi:10.1128/AEM.67.9.4242.

<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=93153&tool=pmcentrez&rendertype=abstract>.

- Reysenbach, Anna-Louise, Yitai Liu, Amy B Banta, Terry J Beveridge, Julie D Kirshtein, Stefan Schouten, Margaret K Tivey, Karen L Von Damm, and Mary a Voytek. 2006. "A Ubiquitous Thermoacidophilic Archaeon from Deep-sea Hydrothermal Vents." *Nature* 442 (7101) (July 27): 444–7. doi:10.1038/nature04921. <http://www.ncbi.nlm.nih.gov/pubmed/16871216>.
- Rinke, Christian, Patrick Schwientek, Alexander Sczyrba, Natalia N. Ivanova, Iain J. Anderson, Jan-Fang Cheng, Aaron Darling, et al. 2013. "Insights into the Phylogeny and Coding Potential of Microbial Dark Matter." *Nature* advance on (July 14): 1–7. doi:10.1038/nature12352. <http://dx.doi.org/10.1038/nature12352>.
- Rivera, Maria C, and James a Lake. 2004. "The Ring of Life Provides Evidence for a Genome Fusion Origin of Eukaryotes." *Nature* 431 (7005) (September 9): 152–5. doi:10.1038/nature02848. <http://www.ncbi.nlm.nih.gov/pubmed/15356622>.
- Roh, Seong Woon, Young-Do Nam, Seong-Hyeuk Nam, Sang-Haeng Choi, Hong-Seog Park, and Jin-Woo Bae. 2010. "Complete Genome Sequence of Halalkalicoccus Jeotgali B3(T), an Extremely Halophilic Archaeon." *Journal of Bacteriology* 192 (17) (September): 4528–9. doi:10.1128/JB.00663-10. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2937367&tool=pmcentrez&render type=abstract>.
- Ronquist, Fredrik, Maxim Teslenko, Paul van der Mark, Daniel L Ayres, Aaron Darling, Sebastian Höhna, Bret Larget, Liang Liu, Marc A Suchard, and John P Huelsenbeck. 2012. "MrBayes 3.2: Efficient Bayesian Phylogenetic Inference and Model Choice Across a Large Model Space." *Systematic Biology* 61 (3) (May): 539–42. doi:10.1093/sysbio/sys029. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3329765&tool=pmcentrez&render type=abstract>.
- Rudolph, Christian, Gerhard Wanner, and Robert Huber. 2001. "Natural Communities of Novel Archaea and Bacteria Growing in Cold Sulfurous Springs with a String-of-Pearls-Like Morphology" 67 (5): 2336–2344. doi:10.1128/AEM.67.5.2336.
- Saitou, N, and M Nei. 1987. "The Neighbor-joining Method: a New Method for Reconstructing Phylogenetic Trees." *Molecular Biology and Evolution* 4 (4) (July): 406–25. <http://www.ncbi.nlm.nih.gov/pubmed/3447015>.
- Sakai, Sanae, Hiroyuki Imachi, Satoshi Hanada, Akiyoshi Ohashi, Hideki Harada, and Yoichi Kamagata. 2008. "Methanocella Paludicola Gen. Nov., Sp. Nov., a Methane-producing Archaeon, the First Isolate of the Lineage 'Rice Cluster I', and Proposal of the New Archaeal Order Methanocellales Ord. Nov." *International Journal of Systematic and Evolutionary Microbiology* 58 (Pt 4) (April): 929–36. doi:10.1099/ijs.0.65571-0. <http://www.ncbi.nlm.nih.gov/pubmed/18398197>.
- Schleper, C, G Puehler, I Holz, a Gambacorta, D Janekovic, U Santarius, H P Klenk, and W Zillig. 1995. "Picrophilus Gen. Nov., Fam. Nov.: a Novel Aerobic, Heterotrophic, Thermoacidophilic Genus and Family Comprising Archaea Capable of Growth Around pH 0." *Journal of Bacteriology* 177 (24) (December): 7050–9. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=177581&tool=pmcentrez&render type=abstract>.

- Schleper, Christa, German Jurgens, and Melanie Jonuscheit. 2005. "Genomic Studies of Uncultivated Archaea." *Nature Reviews. Microbiology* 3 (6) (June): 479–88. doi:10.1038/nrmicro1159. <http://www.ncbi.nlm.nih.gov/pubmed/15931166>.
- Shi, Yanmei, Gene W Tyson, and Edward F DeLong. 2009. "Metatranscriptomics Reveals Unique Microbial Small RNAs in the Ocean's Water Column." *Nature* 459 (7244) (May 14): 266–9. doi:10.1038/nature08055. <http://www.ncbi.nlm.nih.gov/pubmed/19444216>.
- Siebers, Bettina, Melanie Zaparty, Guenter Raddatz, Britta Tjaden, Sonja-Verena Albers, Steve D Bell, Fabian Blombach, et al. 2011. "The Complete Genome Sequence of Thermoproteus Tenax: a Physiologically Versatile Member of the Crenarchaeota." *PLoS One* 6 (10) (January): e24222. doi:10.1371/journal.pone.0024222. http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3189178&tool=pmcentrez&render_type=abstract.
- Slesarev, Alexei I, Katja V Mezhevaya, Kira S Makarova, Nikolai N Polushin, Olga V Shcherbinina, Vera V Shakhova, Galina I Belova, et al. 2002. "The Complete Genome of Hyperthermophile Methanopyrus Kandleri AV19 and Monophyly of Archaeal Methanogens." *Proceedings of the National Academy of Sciences of the United States of America* 99 (7) (April 2): 4644–9. doi:10.1073/pnas.032671499. http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=123701&tool=pmcentrez&render_type=abstract.
- Sogin, S J, M L Sogin, and C R Woese. 1972. "Phylogenetic Measurement in Prokaryotes by Primary Structural Characterization." *Journal of Molecular Evolution* 1 (1) (January): 173–84. <http://www.ncbi.nlm.nih.gov/pubmed/5006250>.
- Spring, Stefan, Reinhard Rachel, Alla Lapidus, Karen Davenport, Hope Tice, Alex Copeland, Jan-Fang Cheng, et al. 2010. "Complete Genome Sequence of Thermosphaera Aggregans Type Strain (M11TL)." *Standards in Genomic Sciences* 2 (3) (January): 245–59. doi:10.4056/sigs.821804. http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3035292&tool=pmcentrez&render_type=abstract.
- Stahl, David a, and José R de la Torre. 2012. "Physiology and Diversity of Ammonia-oxidizing Archaea." *Annual Review of Microbiology* 66 (January): 83–101. doi:10.1146/annurev-micro-092611-150128. <http://www.ncbi.nlm.nih.gov/pubmed/22994489>.
- Stamatakis, Alexandros. 2006. "RAxML-VI-HPC: Maximum Likelihood-based Phylogenetic Analyses with Thousands of Taxa and Mixed Models." *Bioinformatics (Oxford, England)* 22 (21) (November 1): 2688–90. doi:10.1093/bioinformatics/btl446. <http://www.ncbi.nlm.nih.gov/pubmed/16928733>.
- Stamatakis, Alexandros, Paul Hoover, and Jacques Rougemont. 2008. "A Rapid Bootstrap Algorithm for the RAxML Web Servers." *Systematic Biology* 57 (5) (October 1): 758–71. doi:10.1080/10635150802429642. <http://sysbio.oxfordjournals.org/content/57/5/758.short>.
- Stanier, R. Y., Doudoroff, M. & Adelberg, E. A. 1963. *The Microbial World 2nd Edn*. Prentice-Hall, Englewood Cliffs.
- Stein, J L, T L Marsh, K Y Wu, H Shizuya, and E F DeLong. 1996. "Characterization of Uncultivated Prokaryotes: Isolation and Analysis of a 40-kilobase-pair Genome Fragment from a Planktonic Marine Archaeon." *Journal of Bacteriology* 178 (3) (February): 591–9.

<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=177699&tool=pmcentrez&rendertype=abstract>.

- Susanti, Dwi, Eric F Johnson, Jason R Rodriguez, Iain Anderson, Anna a Perevalova, Nikos Kyrpides, Susan Lucas, et al. 2012. "Complete Genome Sequence of *Desulfurococcus Fermentans*, a Hyperthermophilic Cellulolytic Crenarchaeon Isolated from a Freshwater Hot Spring in Kamchatka, Russia." *Journal of Bacteriology* 194 (20) (October): 5703–4. doi:10.1128/JB.01314-12.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3458677&tool=pmcentrez&rendertype=abstract>.
- Takai, K, and K Horikoshi. 1999. "Genetic Diversity of Archaea in Deep-sea Hydrothermal Vent Environments." *Genetics* 152 (4) (August): 1285–97.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1460697&tool=pmcentrez&rendertype=abstract>.
- Tatusov, R L, E V Koonin, and D J Lipman. 1997. "A Genomic Perspective on Protein Families." *Science (New York, N.Y.)* 278 (5338) (October 24): 631–7.
<http://www.ncbi.nlm.nih.gov/pubmed/9381173>.
- Tatusov, R L, D a Natale, I V Garkavtsev, T a Tatusova, U T Shankavaram, B S Rao, B Kiryutin, M Y Galperin, N D Fedorova, and E V Koonin. 2001. "The COG Database: New Developments in Phylogenetic Classification of Proteins from Complete Genomes." *Nucleic Acids Research* 29 (1) (January 1): 22–8.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=29819&tool=pmcentrez&rendertype=abstract>.
- Tatusov, Roman L, Natalie D Fedorova, John D Jackson, Aviva R Jacobs, Boris Kiryutin, Eugene V Koonin, Dmitri M Krylov, et al. 2003. "The COG Database: An Updated Version Includes Eukaryotes." *BMC Bioinformatics* 4 (September 11): 41. doi:10.1186/1471-2105-4-41.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=222959&tool=pmcentrez&rendertype=abstract>.
- Tourna, Maria, Michaela Stieglmeier, Anja Spang, Martin Könneke, Arno Schintlmeister, Tim Urich, Marion Engel, et al. 2011. "Nitrososphaera Viennensis, an Ammonia Oxidizing Archaeon from Soil." *Proceedings of the National Academy of Sciences of the United States of America* 108 (20) (May 17): 8420–5. doi:10.1073/pnas.1013488108.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3100973&tool=pmcentrez&rendertype=abstract>.
- Treusch, Alexander H, Sven Leininger, Arnulf Kletzin, Stephan C Schuster, Hans-Peter Klenk, and Christa Schleper. 2005. "Novel Genes for Nitrite Reductase and Amo-related Proteins Indicate a Role of Uncultivated Mesophilic Crenarchaeota in Nitrogen Cycling." *Environmental Microbiology* 7 (12) (December): 1985–95. doi:10.1111/j.1462-2920.2005.00906.x.
<http://www.ncbi.nlm.nih.gov/pubmed/16309395>.
- Von Jan, Mathias, Alla Lapidus, Tijana Glavina Del Rio, Alex Copeland, Hope Tice, Jan-Fang Cheng, Susan Lucas, et al. 2010. "Complete Genome Sequence of *Archaeoglobus Profundus* Type Strain (AV18)." *Standards in Genomic Sciences* 2 (3) (January): 327–46. doi:10.4056/sigs.942153.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3035285&tool=pmcentrez&rendertype=abstract>.

- Walker, C B, J R de la Torre, M G Klotz, H Urakawa, N Pinel, D J Arp, C Brochier-Armanet, et al. 2010. "Nitrosopumilus Maritimus Genome Reveals Unique Mechanisms for Nitrification and Autotrophy in Globally Distributed Marine Crenarchaea." *Proceedings of the National Academy of Sciences of the United States of America* 107 (19) (May 11): 8818–23. doi:10.1073/pnas.0913533107. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2889351&tool=pmcentrez&render type=abstract>.
- Waters, Elizabeth, Michael J Hohn, Ivan Ahel, David E Graham, Mark D Adams, Mary Barnstead, Karen Y Beeson, et al. 2003. "The Genome of Nanoarchaeum Equitans: Insights into Early Archaeal Evolution and Derived Parasitism." *Proceedings of the National Academy of Sciences of the United States of America* 100 (22) (October 28): 12984–8. doi:10.1073/pnas.1735403100. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=240731&tool=pmcentrez&render type=abstract>.
- Wheelis, M L, O Kandler, and C R Woese. 1992. "On the Nature of Global Classification." *Proceedings of the National Academy of Sciences of the United States of America* 89 (April): 2930–4. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=48777&tool=pmcentrez&render type=abstract>.
- Whittaker, R H. 1969. "New Concepts of Kingdoms or Organisms. Evolutionary Relations Are Better Represented by New Classifications Than by the Traditional Two Kingdoms." *Science (New York, N.Y.)* 163 (3863) (January 10): 150–60. <http://www.ncbi.nlm.nih.gov/pubmed/5762760>.
- Williams, Tom a, Peter G Foster, Tom M W Nye, Cymon J Cox, and T Martin Embley. 2012. "A Congruent Phylogenomic Signal Places Eukaryotes Within the Archaea." *Proceedings. Biological Sciences / The Royal Society* 279 (1749) (December 22): 4870–9. doi:10.1098/rspb.2012.1795. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3497233&tool=pmcentrez&render type=abstract>.
- Woese, C R, and G E Fox. 1977. "Phylogenetic Structure of the Prokaryotic Domain: The Primary Kingdoms." *Proceedings of the National Academy of Sciences of the United States of America* 74 (11) (November): 5088–90. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=432104&tool=pmcentrez&render type=abstract>.
- Woese, C R, O Kandler, and M L Wheelis. 1990. "Towards a Natural System of Organisms: Proposal for the Domains Archaea, Bacteria, and Eucarya." *Proceedings of the National Academy of Sciences of the United States of America* 87 (12) (June): 4576–9. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=54159&tool=pmcentrez&render type=abstract>.
- Woese, C R, L J Magrum, and G E Fox. 1978. "Archaeobacteria." *Journal of Molecular Evolution* 11 (3) (August 2): 245–51. <http://www.ncbi.nlm.nih.gov/pubmed/691075>.
- Woese, C R, J Maniloff, and L B Zablen. 1980. "Phylogenetic Analysis of the Mycoplasmas." *Proceedings of the National Academy of Sciences of the United States of America* 77 (1) (January): 494–8.

<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=348298&tool=pmcentrez&rendertype=abstract>.

Woese, C R, G J Olsen, M Ibba, and D Söll. 2000. "Aminoacyl-tRNA Synthetases, the Genetic Code, and the Evolutionary Process." *Microbiology and Molecular Biology Reviews* : *MMBR* 64 (1) (March): 202–36.

<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=98992&tool=pmcentrez&rendertype=abstract>.

Woese, Carl R, Kenneth R Luehrsen, Cheryl D Pribula, and George E Fox. 1976. "Molecular Evolution Sequence Characterization of 5S Ribosomal RNA from Eight Gram Positive Procaryotes" 8: 143–153.

Woyke, Tanja, Damon Tighe, Konstantinos Mavromatis, Alicia Clum, Alex Copeland, Wendy Schackwitz, Alla Lapidus, et al. 2010. "One Bacterial Cell, One Complete Genome." *PLoS One* 5 (4) (January): e10314. doi:10.1371/journal.pone.0010314.

<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2859065&tool=pmcentrez&rendertype=abstract>.

Yutin, Natalya, Kira S Makarova, Sergey L Mekhedov, Yuri I Wolf, and Eugene V Koonin. 2008. "The Deep Archaeal Roots of Eukaryotes." *Molecular Biology and Evolution* 25 (8) (August): 1619–30. doi:10.1093/molbev/msn108.

<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2464739&tool=pmcentrez&rendertype=abstract>.

Zillig, W., K.O. Stetter, W. Schäfer, D. Janekovic, S. Wunderl, I. Holz, and P. Palm. 1981. "Thermoproteales: A Novel Type of Extremely Thermoacidophilic Anaerobic Archaeobacteria Isolated from Icelandic Solfataras." *Zentralblatt Für Bakteriologie Mikrobiologie Und Hygiene: I. Abt. Originale C: Allgemeine, Angewandte Und Ökologische Mikrobiologie* 2 (3): 227–205. doi:10.1016/S0721-9571(81)80001-4. [http://dx.doi.org/10.1016/S0721-9571\(81\)80001-4](http://dx.doi.org/10.1016/S0721-9571(81)80001-4).

Zuckerandl, E, and L Pauling. 1965. "Molecules as Documents of Evolutionary History." *Journal of Theoretical Biology* 8 (2) (March): 357–66. <http://www.ncbi.nlm.nih.gov/pubmed/5876245>.

Zvyagintseva, I.S., and A.L. Tarasov. 1988. "Extreme Halophilic Bacteria from Saline Soils." *Microbiology* v. 56(5) p (March 1). <http://agris.fao.org/agris-search/search/display.do?f=1989/US/US89321.xml;US8858510>.

Annexes

Table des Annexes

Annexe 1 : Matériels Supplémentaires de l'Article 1 (Chapitre 1).....	- 2 -
Annexe 2 : Listes et annotation des protéines des trois génomes étudiées.....	- 36 -
a. Liste des protéines de <i>Nitrosopumilus maritimus</i>	- 36 -
b. Liste des protéines de <i>Cenarchaeum symbiosum</i>	- 36 -
c. Liste des protéines de ' <i>Candidatus Caldiarchaeum subterraneum</i> '	- 36 -
Annexe 3 : Répartition taxonomique des 200 nouveaux marqueurs chez les archées.	- 60 -
Annexe 4 : Liste des Protéines informationnelles mises à jour pour le Chapitre 1.	- 66 -
Annexe 5 : Liste des matrices utilisées pour l'analyse en désaturation par sélection de sites pour l'étude exposée dans le Chapitre 1.	- 68 -
Annexe 6 : Liste des matrices utilisées pour l'analyse en désaturation par sélection de gènes pour l'étude exposée dans le Chapitre 1.	- 70 -
Annexe 7 : Matériels Supplémentaires de l'Article 2 (Chapitre 2).....	- 78 -
Annexe 8 : Liste des 38 nouveaux marqueurs sélectionnés pour l'étude présentés dans le Chapitre 2.....	- 108 -
Annexe 9 : Répartition taxonomique des 38 nouveaux marqueurs chez les archées et chez les bactéries.	- 110 -
Annexe 10 : Liste des matrices utilisées pour l'analyse en désaturation par sélection de sites pour l'étude exposée dans le Chapitre 2.	- 114 -
Annexe 11 : Liste des matrices utilisées pour l'analyse en désaturation par sélection de gènes pour l'étude exposée dans le Chapitre 2.	- 116 -
Annexe 12 : Matériels Supplémentaires de l'Article 3 (Chapitre 3).....	- 120 -

Annexe 1 : Matériels supplémentaires de l'Article 1 (Chapitre 1).

Supplementary material.

Supplementary Table S1. Table showing the taxonomic sampling used for the identification of new markers that can be potentially used to study the phylogeny of Archaea.

Supplementary Table S2. Table showing the 200 new proteins used for the phylogenetic analysis.

Supplementary Table S3. Table showing the taxonomic distribution of the 129 archaeal genomes used for the assembly of final datasets.

Supplementary Figure S1. Unrooted maximum likelihood phylogeny of Archaea based on the 200 newly identified proteins (L2 supermatrix, 48,904 amino acid positions, 122 species). The tree was inferred with PhyML (LG+ Γ 8). The scale bar represents the average number of substitutions per site. Numbers at branches are SH-like supports

Supplementary Figure S2. Unrooted maximum likelihood phylogeny of Archaea based on the 53 ribosomal proteins (L3 supermatrix, 6,228 amino acid positions, 122 organisms). The tree was inferred with PhyML (LG+ Γ 8). The scale bar represents the average number of substitutions per site. Numbers at branches are SH-like supports.

Supplementary Figure S3. Unrooted maximum likelihood phylogeny of Archaea based on the 15 subunits of the RNA polymerase and transcription factors (L4 supermatrix, 2,970 amino acid positions, 122 species). The tree was inferred with PhyML (LG+ Γ 8). The scale bar represents the average number of substitutions per site. Numbers at branches are SH-like supports.

Supplementary Figure S4. Unrooted Bayesian tree of Archaea based on the 15 subunits of the RNA polymerase and transcription factors (L4 supermatrix, 2970 amino acid positions, 122 species). The tree was inferred with PhyloBayes (CAT+ Γ 8). The scale bar represents the average number of substitutions per site. Numbers at branches are posterior probabilities.

Supplementary Figure S5. Unrooted maximum likelihood trees resulting from the site-by-site desaturation strategy. The S_0 to S_{34} supermatrices were extracted from the XL1 supermatrix (see material and methods). Trees were inferred with PhyML (CAT+ Γ 8). The scale bar represents the average number of substitutions per site. Numbers at branches are SH supports.

Eucarya

Allomyces macrogynus
Arabidopsis thaliana
Arabidopsis thaliana Plastid
Aspergillus niger
Aureococcus anophagefferens
Aureococcus anophagefferens Plastid
Batrachochytrium dendrobatidis JAM81
Bigelowiella natans Nucleomorph
Bigelowiella natans Plastid
Blastocystis hominis
Bodo saltans
Botrytis cinerea
Brachypodium distachyon
Brachypodium distachyon Plastid
Caenorhabditis elegans
Calliarthron tuberculosum
Capsaspora owczarzaki
Chaetomium globosum
Chlamydomonas reinhardtii
Chlamydomonas reinhardtii Plastid
Chlorella sp. NC64A
Chlorella vulgaris C-169
Chlorella vulgaris Plastid
Chondrus crispus EST
Chromera velia Plastid
Coprinopsis cinerea
Cryptococcus neoformans JEC21
Cryptosporidium parvum
Cyanidioschyzon merolae
Cyanidioschyzon merolae Plastid
Cyanophora paradoxa
Cyanophora paradoxa EST
Cyanophora paradoxa Plastid
Cyanophora paradoxa T
Daphnia pulex
Debaryomyces hansenii CBS767
Dictyostelium discoideum
Dictyostelium purpureum QSDP1
Diplonema papillatum
Ectocarpus siliculosus
Ectocarpus siliculosus Plastid
Emiliana huxleyi CCMP1516
Emiliana huxleyi Plastid
Entamoeba histolytica
Eucheuma denticulatum EST
Euglena gracilis
Fragilariopsis cylindrus
Furcellaria lumbricalis EST
Fusarium graminearum
Galdieria sulphuraria
Giardia lamblia
Gracilaria changii EST
Griffithsia okiensis EST
Guillardia theta Nucleomorph
Guillardia theta Plastid
Helobdella robusta

Histoplasma capsulatum
Homo sapiens
Laccaria bicolor
Leishmania infantum
Leishmania major strain Friedlin
Melampsora larici-populina
Micromonas pusilla CCMP1545
Micromonas pusilla CCMP1545 Plastid
Micromonas sp. RCC299 Plastid
Micromonas strain RCC299
Mimulus guttatus
Monosiga brevicollis
Mucor circinelloides
Mycosphaerella graminicola
Naegleria gruberi
Nematostella vectensis
Neurospora crassa OR74A
Oryza sativa
Oryza sativa Plastid
Ostreococcus lucimarinus
Ostreococcus tauri
Ostreococcus tauri Plastid
Paramecium tetraurelia
Perkinsus marinus
Phaeodactylum tricornutum
Phaeodactylum tricornutum Plastid
Phanerochaete chrisosporium
Phycomyces blakesleeianus
Physcomitrella patens
Physcomitrella patens Plastid
Phytophthora capsici
Phytophthora sojae
Plasmodium falciparum 3D7
Populus trichocarpa
Populus trichocarpa Plastid
Porphyra sp EST
Porphyra yezoensis Plastid
Porphyridium cruentum EST
Postia placenta
Proccryptobia sorokini
Puccinia graminis
Rhizopus oryzae
Rhynchomonas nasuta
Saccharomyces cerevisiae RM11-1a
Salpingoeca rosetta
Schizophyllum commune
Schizosaccharomyces pombe
Selaginella moellendorffii
Selaginella moellendorffii Plastid
Spizellomyces punctatus
Sporobolomyces roseus
Stagonospora nodorum
Tetrahymena thermophila
Thalassiosira pseudonana
Thalassiosira pseudonana Plastid
Thecamonas trahens
Toxoplasma gondii GT1
Toxoplasma gondii Plastid

Trichoderma reesei
Trichomonas vaginalis G3
Trichoplax adhaerens Grell-BS-1999
Trypanosoma brucei TREU927
Trypanosoma cruzi strain CL Brener
Ustilago maydis
Volvox carteri
Yarrowia lipolytica

Bacteria

Acaryochloris marina MBIC11017
Acholeplasma laidlawii PG_8A
Acidiphilium cryptum JF_5
Acidithiobacillus ferrooxidans ATCC 23270
Acidobacteria bacterium Ellin345
Acidothermus cellulolyticus 11B
Acinetobacter baumannii AB0057
Actinobacillus pleuropneumoniae L20
Aeromonas hydrophila ATCC_7966
Agrobacterium tumefaciens C58
Akkermansia muciniphila ATCC BAA-835
Alcanivorax borkumensis SK2
Aliivibrio salmonicida LFI1238
Alkalilimnicola ehrlichei MLHE_1
Alkaliphilus metalliredigens QYMF
Alteromonas macleodii Deep ecotype
Amoebophilus asiaticus 5a2
Anabaena variabilis
Anaerocellum thermophilum DSM 6725
Anaeromyxobacter sp. Fw109_5
Anaplasma marginale StMaries
Anoxybacillus flavithermus WK1
Aquifex aeolicus VF5
Arcobacter butzleri RM4018
Aromatoleum aromaticum EbN1
Arthrobacter aurescens TC1
Aster yellows witches-broom phytoplasma AYWB
Azoarcus sp. BH72
Azorhizobium caulinodans ORS_571
Bacillus amyloliquefaciens FZB42
Bacteroides fragilis NCTC_9343
Bacteroides fragilis YCH46
Bartonella bacilliformis KC583
Baumannia cicadellinicola Hc
Bdellovibrio bacteriovorus HD100
Beijerinckia indica subsp. *indica* ATCC 9039
Bifidobacterium adolescentis ATCC_15703
Blochmannia floridanus
Bordetella bronchiseptica RB50
Borrelia afzelii PKo
Brachyspira hyodysenteriae WA1
Bradyrhizobium sp. BTAi1
Brucella abortus 9_941
Buchnera aphidicola str. APS
Burkholderia sp. 383
Caldicellulosiruptor saccharolyticus DSM_8903

Campylobacter concisus 13826
Carboxydotherrmus hydrogenoformans Z_2901
Carsonella ruddii PV
Caulobacter crescentus CB15
Cellvibrio japonicus Ueda107
Chlamydia muridarum Nigg
Chlamydia trachomatis 70
Chlamydophila abortus S26_3
Chlamydophila pneumoniae CWL029
Chlorobaculum parvum NCIB 8327
Chlorobium chlorochromatii CaD3
Chloroflexus aurantiacus J_10_fl
Chloroherpeton thalassium ATCC 35110
Chromobacterium violaceum ATCC_12472
Chromohalobacter salexigens DSM 3043
Citrobacter koseri ATCC_BAA895
Clavibacter michiganensis NCPPB_382
Clostridium acetobutylicum ATCC_824
Colwellia psychrerythraea 34H
Coprothermobacter proteolyticus DSM 5265
Corynebacterium diphtheriae NCTC_13129
Coxiella burnetii RSA_493
Crocospaera watsonii WH8501
Cronobacter sakazakii ATCC BAA-894
Cupriavidus taiwanensis
Cyanobacterium Yellowstone A
Cyanobacterium Yellowstone B
Cyanothece sp. ATCC51142
Cyanothece sp. CCY0110
Cytophaga hutchinsonii ATCC_33406
Dechloromonas aromatica RCB
Dehalococcoides sp. BAV1
Deinococcus geothermalis DSM_11300
Delftia acidovorans SPH_1
Desulfitobacterium hafniense Y51
Desulfococcus oleovorans Hxd3
Desulfotalea psychrophila LSv54
Desulfotomaculum reducens MI_1
Desulfovibrio desulfuricans G20
Desulfovibrio vulgaris DP4
Desulfovibrio vulgaris Hildenborough
Diaphorobacter sp. TPSY
Dichelobacter nodosus VCS1703A
Dictyoglomus thermophilum H-6-12
Dinoroseobacter shibae DFL_12
Ehrlichia canis Jake
Elusimicrobium minutum Pei191
Enterobacter sp. 638
Enterococcus faecalis V583
Erythrobacter litoralis HTCC2594
Escherichia coli 536
Exiguobacterium sibiricum 255-15
Fervidobacterium nodosum Rt17_B1
Finegoldia magna ATCC_29328
Flavobacterium johnsoniae UW101
Flavobacterium psychrophilum JIP02_86
Francisella novicida U112
Frankia alni ACN14a

Fusobacterium nucleatum ATCC_25586
Geobacillus kaustophilus HTA426
Geobacter metallireducens GS_15
Geobacter sulfurreducens PCA
Geobacter uraniumreducens Rf4
Gloeobacter violaceus PCC7421
Gluconacetobacter diazotrophicus PAI_5
Gramella forsetii KT0803
Granulibacter bethesdensis CGDNIH1
Haemophilus ducreyi 35000HP
Hahella chejuensis KCTC_2396
Haliangium ochraceum DSM 14365
Halorhodospira halophila SL1
Halothermothrix orenii H 168
Helicobacter acinonychis Sheeba
Heliobacterium modesticaldum Ice1
Hermiimonas arsenicoxydans
Herpetosiphon aurantiacus ATCC_23779
Hydrogenobaculum sp. Y04AAS1
Hyphomonas neptunium ATCC_15444
Idiomarina loihiensis L2TR
Jannaschia sp. CCS1
Janthinobacterium sp. Marseille
Kineococcus radiotolerans SRS30216
Klebsiella pneumoniae MGH_78578
Kocuria rhizophila DC2201
Lactobacillus acidophilus NCFM
Lactococcus lactis II1403
Lawsonia intracellularis PHE_MN1_00
Legionella pneumophila Corby
Leifsonia xyli CTCB07
Leptospira borgpetersenii JB197
Leptothrix cholodnii SP-6
Leuconostoc citreum KM20
Listeria innocua Clip11262
Lyngbya sp. PCC8106
Lysinibacillus sphaericus C3_41
Macrococcus caseolyticus JCSC5402
Magnetococcus sp. MC_1
Magnetospirillum magneticum AMB_1
Mannheimia succiniciproducens MBEL55E
Maricaulis maris MCS10
Marinobacter aquaeolei VT8
Marinomonas sp. MWYL1
Mesoplasma florum L1
Mesorhizobium sp. BNC1
Methylbium petroleiphilum PM1
Methylobacillus flagellatus KT
Methylobacterium extorquens PA1
Methylocella silvestris BL2
Methylococcus capsulatus Bath
Microcystis aeruginosa NIES843
Moorella thermoacetica ATCC_39073
Mycobacterium abscessus
Mycoplasma agalactiae PG2
Myxococcus xanthus DK_1622
Natranaerobius thermophilus JWNM-WN-LF
Nautilia profundicola AmH

Neisseria gonorrhoeae FA_1090
Neorickettsia sennetsu Miyayama
Nitratiruptor sp. SB155_2
Nitrobacter hamburgensis X14
Nitrosococcus oceani ATCC_19707
Nitrosomonas europaea ATCC_19718
Nitrospira multififormis ATCC_25196
Nocardia farcinica IFM_10152
Nocardioides sp. JS614
Nodularia spumigena CCY9414
Nostoc punctiforme
Nostoc sp. PCC7120
Novosphingobium aromaticivorans DSM_12444
Oceanobacillus iheyensis HTE831
Ochrobactrum anthropi ATCC_49188
Oenococcus oeni PSU_1
Oligotropha carboxidovorans OM5
Onion yellows phytoplasma OY-M
Opitutus terrae PB90-1
Parabacteroides distasonis ATCC_8503
Parachlamydia acanthamoebae
Paracoccus denitrificans PD1222
Parvibaculum lavamentivorans DS1
Pasteurella multocida Pm70
Pectobacterium atrosepticum SCRI1043
Pediococcus pentosaceus ATCC_25745
Pelagibacter ubique HTCC1062
Pelobacter carbinolicus DSM_2380
Pelobacter propionicus DSM_2379
Pelodictyon luteolum DSM_273
Pelotomaculum thermopropionicum SI
Petrotoga mobilis SJ95
Phenylobacterium zucineum HLK1
Photobacterium profundum SS9
Photorhabdus luminescens TTO1
Phytoplasma australiense
Plesiocystis pacifica SIR-1
Polaromonas sp. JS666
Polynucleobacter sp. QLW_P1DMWA_1
Porphyromonas gingivalis W83
Prochlorococcus marinus AS9601
Prochlorococcus marinus CCMP1375
Prochlorococcus marinus CCMP1986
Prochlorococcus marinus MIT9215
Prochlorococcus marinus MIT9301
Prochlorococcus marinus MIT9303
Prochlorococcus marinus MIT9312
Prochlorococcus marinus MIT9313
Prochlorococcus marinus MIT9515
Prochlorococcus marinus NATL1A
Prochlorococcus marinus NATL2A
Propionibacterium acnes KPA171202
Prosthecochloris vibrioformis DSM_265
Proteus mirabilis HI4320
Protochlamydia amoebophila UWE25
Pseudoalteromonas atlantica T6c
Pseudomonas aeruginosa PAO1
Psychrobacter arcticus 273_4

Psychromonas ingrahamii 37
Ralstonia eutropha H16
Renibacterium salmoninarum ATCC_33209
Rhizobium etli CFN_42
Rhodobacter sphaeroides 2_4_1
Rhodococcus sp. RHA1
Rhodoferax ferrireducens T118
Rhodopirellula baltica SH_1
Rhodopseudomonas palustris BisA53
Rhodospirillum rubrum ATCC_11170
Rickettsia akari Hartford
Roseiflexus castenholzii DSM_13941
Roseobacter denitrificans OCh_114
Rubrobacter xylanophilus DSM_9941
Ruegeria sp. TM1040
Saccharophagus degradans 2_40
Saccharopolyspora erythraea NRRL_2338
Salinibacter ruber DSM_13855
Salinispora arenicola CNS_205
Salmonella enterica arizonae
Serratia proteamaculans 568
Shewanella amazonensis SB2B
Shigella boydii Sb227
Silicibacter pomeroyi DSS_3
Sinorhizobium medicae WSM419
Sodalis glossinidius
Solibacter usitatus Ellin6076
Sorangium cellulosum
Sphingomonas wittichii RW1
Sphingopyxis alaskensis RB2256
Staphylococcus aureus MRSA252
Stenotrophomonas maltophilia K279a
Streptococcus agalactiae 2603V_R
Streptomyces avermitilis MA_4680
Sulcia muelleri GWSS
Sulfurihydrogenibium sp. YO3AOP1
Sulfurimonas denitrificans DSM 1251
Sulfurovum sp. NBC37_1
Symbiobacterium thermophilum IAM_14863
Synechococcus elongatus PCC6301
Synechococcus elongatus PCC7942
Synechococcus sp. WH5701
Synechococcus sp. CC9311
Synechococcus sp. CC9605
Synechococcus sp. CC9902
Synechococcus sp. RCC307
Synechococcus sp. WH7803
Synechococcus sp. WH8102
Synechocystis sp. PCC6803
Syntrophobacter fumaroxidans MPOB
Syntrophomonas wolfei subsp. *wolfei* str. Goettingen
Syntrophus aciditrophicus SB
Thermoanaerobacter pseudethanolicus ATCC_33223
Thermobifida fusca YX
Thermodesulfovibrio yellowstonii DSM 11347
Thermomicrobium roseum DSM 5159
Thermosiphon melanesiensis BI429
Thermosynechococcus elongatus BP1

Thermotoga lettingae TMO
Thermus thermophilus HB8
Thioalkalivibrio sp. HL-EbGR7
Thiobacillus denitrificans ATCC_25259
Thiomicrospira crunogena XCL_2
Treponema denticola ATCC_35405
Trichodesmium erythraeum
Tropheryma whipplei TW08_27
Ureaplasma parvum ATCC_700970
Verminephrobacter eiseniae EF01_2
Vibrio cholerae N16961
Wigglesworthia glossinidia
Wolbachia endosymbiont TRS
Wolinella succinogenes DSM_1740
Xanthobacter autotrophicus Py2
Xanthomonas campestris ATCC_33913
Xylella fastidiosa 9a5c
Yersinia enterocolitica 8081
Zymomonas mobilis ZM4

Archaea

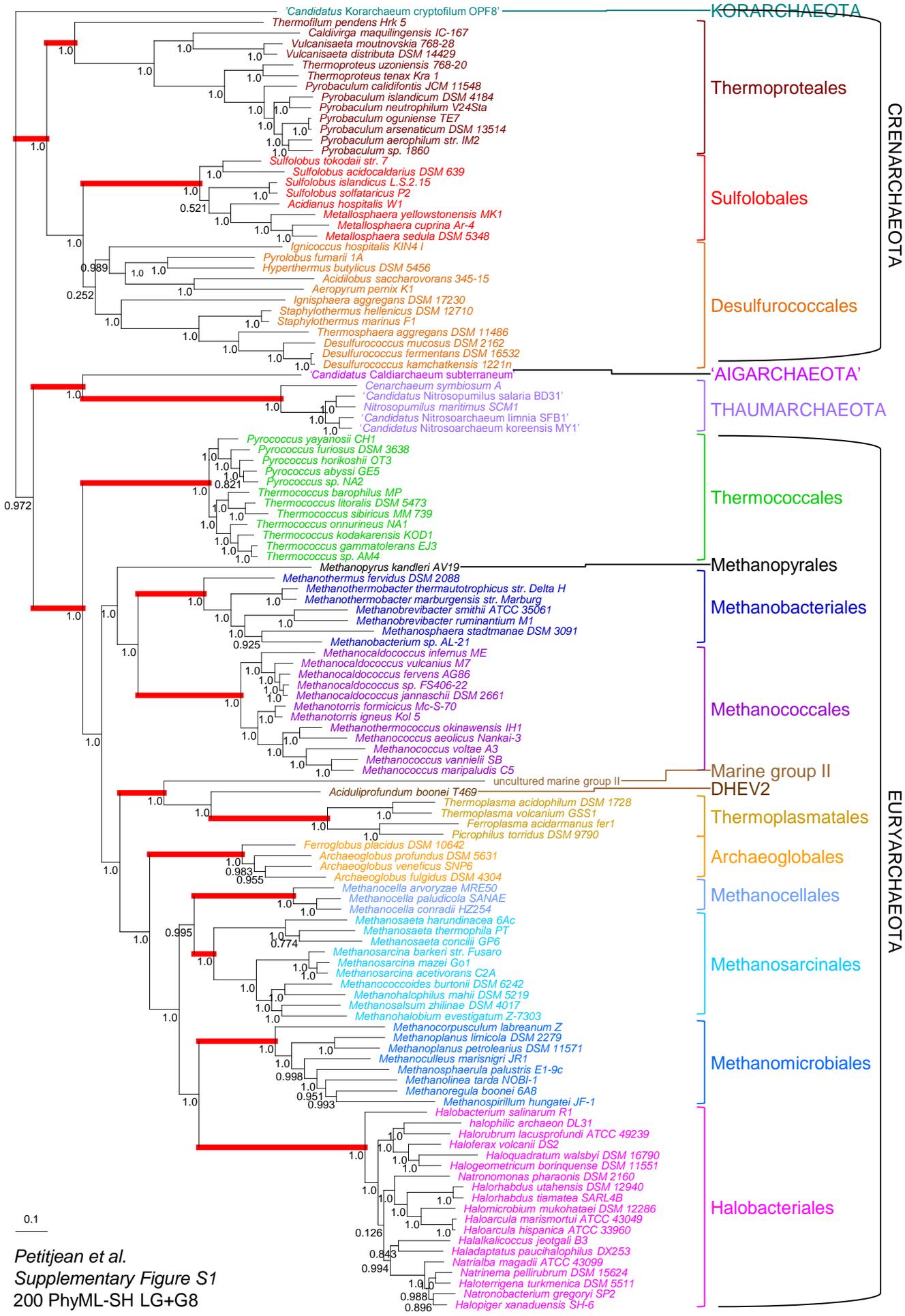
Acidilobus saccharovorans 345-15
Aciduliprofundum boonei T469
Aeropyrum pernix K1
Archaeoglobus fulgidus DSM_4304
Archaeoglobus profundus DSM 5631
Caldivirga maquilingensis IC_167
Candidatus Caldiarchaeum subterraneum
Candidatus Korarchaeum cryptofilum OPF8
Candidatus Methanoregula boonei 6A8
Candidatus Micrarchaeum acidiphilum ARMAN-2
Candidatus Parvarchaeum acidiphilum ARMAN-4
Candidatus Parvarchaeum acidophilus ARMAN-5
Cenarchaeum symbiosum A
Desulfurococcus kamchatkensis
Desulfurococcus mucosus DSM2162
Ferroglobus placidus DSM 10642
Ferroplasma acidarmanus fer1
Halalkalicoccus jeotgali B3
Haloarcula marismortui ATCC_43049
Halobacterium salinarum R1
Haloferax volcanii DS2
Halogeometricum borinquense DSM 11551
Halomicrobium mukohataei DSM 12286
Haloquadratum walsbyi DSM_16790
Halorhabdus utahensis DSM 12940
Halorubrum lacusprofundi ATCC 49239
Haloterrigena turkmenica DSM 5511
Hyperthermus butylicus DSM_5456
Ignicoccus hospitalis KIN4_I
Ignisphaera aggregans DSM17230
Metallosphaera sedula DSM_5348
Methanobrevibacter ruminantium M1
Methanobrevibacter smithii ATCC_35061
Methanocaldococcus fervens AG86
Methanocaldococcus infernus ME

Methanocaldococcus jannaschii DSM 2661
Methanocaldococcus sp. FS406-22
Methanocaldococcus vulcanius M7
Methanocella paludicola SANA E
Methanococcoides burtonii DSM_6242
Methanococcus aeolicus Nankai_3
Methanococcus maripaludis C5
Methanococcus vanniellii
Methanococcus voltae A3
Methanocorpusculum labreanum Z
Methanoculleus marisnigri JR1
Methanohalobium evestigatum Z-7303
Methanohalophilus mahii DSM 5219
Methanoplanus petrolearius DSM11571
Methanopyrus kandleri AV19
Methanosaeta thermophila PT
Methanosarcina acetivorans C2A
Methanosarcina barkeri str. Fusaro
Methanosarcina mazei Go1
Methanosphaera stadtmanae DSM_3091
Methanosphaerula palustris E1-9c
Methanospirillum hungatei JF_1
Methanothermobacter marburgensis str. Marburg
Methanothermobacter thermautotrophicus
Methanothermococcus okinawensis IH1
Methanothermus fervidus
Nanoarchaeum equitans Kin4_M
Natrialba magadii ATCC 43099
Natronomonas pharaonis DSM_2160
Nitrosopumilus maritimus SCM1
Picrophilus torridus DSM_9790
Pyrobaculum aerophilum IM2
Pyrobaculum arsenaticum DSM 13514
Pyrobaculum calidifontis JCM 11548
Pyrobaculum islandicum DSM 4184
Pyrobaculum neutrophilum V24Sta
Pyrococcus abyssi GE5
Pyrococcus furiosus DSM 3638
Pyrococcus horikoshii OT3
Staphylothermus hellenicus DSM 12710
Staphylothermus marinus F1
Sulfolobus acidocaldarius DSM_639
Sulfolobus islandicus L.S.2.15
Sulfolobus solfataricus P2
Sulfolobus tokodaii str. 7
Thermococcus barophilus MP
Thermococcus gammatolerans EJ3
Thermococcus kodakarensis KOD1
Thermococcus onnurineus NA1
Thermococcus sibiricus MM 739
Thermococcus sp. AM4
Thermofilum pendens Hrk_5
Thermoplasma acidophilum DSM_1728
Thermoplasma volcanium GSS1
Thermosphaera aggregans DSM 11486
uncultured methanogenic archaeon RC-1
Vulcanisaeta distributa DSM14429

155	hypothetical protein Nmar_1233	161528741	YP_001582567.1	M217	COG1909	S	Function unknown	POORLY CHARACTERIZED	10	122
156	rRNA pseudouridine synthase D TruD	161528706	YP_001582532.1	M218	COG0585	S	Function unknown	POORLY CHARACTERIZED	10	125
157	peptidyl-rRNA hydrolase	161528708	YP_001582534.1	M219	COG1990	S	Function unknown	POORLY CHARACTERIZED	10	127
158	conserved hypothetical protein	315425963	BAJ47612.1	MA25	COG1371	S	Function unknown	POORLY CHARACTERIZED	10	124
159	DNA polymerase II small subunit	161527514	YP_001581340.1	M199	COG4304	S	Function unknown	POORLY CHARACTERIZED	Totale	94
160	GTP cyclohydrolase IIA	161529122	YP_001582948.1	M200	COG2429	S	Function unknown	POORLY CHARACTERIZED	Totale	67
161	hypothetical protein Nmar_0307	161527815	YP_001581641.1	M201	COG1698	S	Function unknown	POORLY CHARACTERIZED	Totale	98
162	hypothetical protein Nmar_0387	161527895	YP_001581721.1	M202	COG1650	S	Function unknown	POORLY CHARACTERIZED	Totale	112
163	hypothetical protein Nmar_0599	161528107	YP_001581933.1	M203	COG4046	S	Function unknown	POORLY CHARACTERIZED	Totale	45
164	hypothetical protein Nmar_0810	161528318	YP_001582144.1	M204	COG2106	S	Function unknown	POORLY CHARACTERIZED	Totale	87
165	hypothetical protein Nmar_0812	161528320	YP_001582146.1	M205	COG2090	S	Function unknown	POORLY CHARACTERIZED	Totale	101
166	hypothetical protein Nmar_1415	161528923	YP_001582749.1	M206	COG2237	S	Function unknown	POORLY CHARACTERIZED	Totale	110
167	hypothetical protein Nmar_1719	161529227	YP_001583053.1	M207	COG1701	S	Function unknown	POORLY CHARACTERIZED	Totale	113
168	LPPG:FO 2-phospho-L-lactate transferase	161528134	YP_001581960.1	M208	COG0391	S	Function unknown	POORLY CHARACTERIZED	Totale	70
169	hypothetical protein Nmar_0053	161527565	YP_001581391.1	M209	COG1679	S	Function unknown	POORLY CHARACTERIZED	Totale	90
170	hypothetical protein Nmar_0054	161527566	YP_001581392.1	M210	COG1786	S	Function unknown	POORLY CHARACTERIZED	Totale	90
171	hypothetical protein Nmar_0617	161528125	YP_001581951.1	M213	COG2042	S	Function unknown	POORLY CHARACTERIZED	Totale	91
172	hypothetical protein Nmar_1606	161529114	YP_001582940.1	M215	COG1303	S	Function unknown	POORLY CHARACTERIZED	Totale	111
173	hypothetical conserved protein	315425224	BAJ46893.1	MA09	COG1839	S	Function unknown	POORLY CHARACTERIZED	Totale	61
174	ATP phosphoribosyltransferase	315425361	BAJ47027.1	MA11	COG4077	S	Function unknown	POORLY CHARACTERIZED	Totale	81
175	conserved hypothetical protein	315426248	BAJ47890.1	MA31	COG1469	S	Function unknown	POORLY CHARACTERIZED	Totale	90
176	ATPase of the PP-loop superfamily	118575616	YP_875359.1	M181	COG2102	R	General function prediction only	POORLY CHARACTERIZED	10	127
177	DNA-binding TFAR19-related protein	161529131	YP_001582957.1	M184	COG2118	R	General function prediction only	POORLY CHARACTERIZED	10	129
178	eRF1 domain-containing 1 protein	161527621	YP_001581447.1	M185	COG1537	R	General function prediction only	POORLY CHARACTERIZED	10	128
179	hypothetical protein Nmar_1562	161529070	YP_001582896.1	M186	COG1571	R	General function prediction only	POORLY CHARACTERIZED	10	120
180	putative RNA-processing protein	161528917	YP_001582743.1	M188	COG1094	R	General function prediction only	POORLY CHARACTERIZED	10	126
181	beta-lactamase domain-containing protein	161529235	YP_001583061.1	M190	COG1782	R	General function prediction only	POORLY CHARACTERIZED	10	129
182	phosphopyruvate hydratase	161527825	YP_001581651.1	M191	COG2102	R	General function prediction only	POORLY CHARACTERIZED	10	121
183	ATPase	161528562	YP_001582388.1	M193	COG1855	R	General function prediction only	POORLY CHARACTERIZED	10	126
184	hypothetical protein Nmar_1571	161529079	YP_001582905.1	M194	COG1634	R	General function prediction only	POORLY CHARACTERIZED	10	119
185	putative ATPase RIL	161527803	YP_001581629.1	M198	COG1245	R	General function prediction only	POORLY CHARACTERIZED	10	123
186	ribonuclease Z	315426519	BAJ48150.1	MA36	COG1234	R	General function prediction only	POORLY CHARACTERIZED	10	120
187	GTP-binding protein	315426744	BAJ48368.1	MA41	COG1163	R	General function prediction only	POORLY CHARACTERIZED	10	127
188	5-formaminoimidazole-4-carboxamide-1-(beta)-D-ribofuranosyl 5p-monophosphate synth	161528278	YP_001582104.1	M182	COG1759	R	General function prediction only	POORLY CHARACTERIZED	Totale	59
189	DEAD/DEAH box helicase domain-containing protein	161527751	YP_001581577.1	M183	COG1204	R	General function prediction only	POORLY CHARACTERIZED	Totale	112
190	nucleotide binding protein PINc	161527947	YP_001581773.1	M187	COG1439	R	General function prediction only	POORLY CHARACTERIZED	Totale	118
191	5-formaminoimidazole-4-carboxamide-1-(beta)-D-ribofuranosyl 5p-monophosphate synth	161528380	YP_001582206.1	M189	COG1759	R	General function prediction only	POORLY CHARACTERIZED	Totale	82
192	SPP-like hydrolase	161528193	YP_001582019.1	M192	COG0561	R	General function prediction only	POORLY CHARACTERIZED	Totale	106
193	GTP1/OBG protein	161529113	YP_001582939.1	M195	COG2262	R	General function prediction only	POORLY CHARACTERIZED	Totale	93
194	GHMP kinase	161529226	YP_001583052.1	M196	COG1829	R	General function prediction only	POORLY CHARACTERIZED	Totale	115
195	geranylgeranylglyceryl phosphate synthase-like protein	161529303	YP_001583129.1	M197	COG1646	R	General function prediction only	POORLY CHARACTERIZED	Totale	87
196	acetylaclyl transferase related protein	315425161	BAJ46831.1	MA06	COG0110	R	General function prediction only	POORLY CHARACTERIZED	Totale	49
197	phosphopantetheine adenylyltransferase	315426113	BAJ47758.1	MA28	COG1019	R	General function prediction only	POORLY CHARACTERIZED	Totale	99
198	nucleic-acid-binding protein	315426921	BAJ48540.1	MA46	COG1545	R	General function prediction only	POORLY CHARACTERIZED	Totale	96
199	conserved hypothetical protein	315427069	BAJ48685.1	MA50	COG1938	R	General function prediction only	POORLY CHARACTERIZED	Totale	107
200	conserved hypothetical protein	315427074	BAJ48690.1	MA52	COG2129	R	General function prediction only	POORLY CHARACTERIZED	Totale	82

TaxID	Species	Phylum	Class/Order
* : species retained for the construction of the final datasets (108 species)			
1	311458 <i>Candidatus</i> Caldiarchaeum subterraneum	* Aigarchaeota	Unclassified
2	414004 <i>Cenarchaeum symbiosum</i> A	* Thaumarchaeota	Cenarchaeales
3	436308 <i>Nitrosopumilus maritimus</i> SCM1	* Thaumarchaeota	Nitrosopumilales
4	886738 <i>Candidatus</i> Nitrosoarchaeum limnia SFB1	* Thaumarchaeota	Nitrosopumilales
5	1001994 <i>Candidatus</i> Nitrosoarchaeum koreensis MY1	* Thaumarchaeota	Nitrosopumilales
6	859350 <i>Candidatus</i> Nitrosopumilus salaria BD31	* Thaumarchaeota	Nitrosopumilales
7	797209 <i>Haladaptatus paucihalophilus</i> DX253	* Euryarchaeota	Halobacteriales
8	795797 <i>Halalkalicoccus jeotgali</i> B3	Euryarchaeota	Halobacteriales
9	634497 <i>Haloarcula hispanica</i> ATCC 33960	Euryarchaeota	Halobacteriales
10	272569 <i>Haloarcula marismortui</i> ATCC 43049	Euryarchaeota	Halobacteriales
11	478009 <i>Halobacterium salinarum</i> R1	Euryarchaeota	Halobacteriales
12	309800 <i>Haloferax volcanii</i> DS2	Euryarchaeota	Halobacteriales
13	469382 <i>Halogeometricum borinquense</i> DSM 11551	* Euryarchaeota	Halobacteriales
14	485914 <i>Halomicrobium mukohataei</i> DSM 12286	Euryarchaeota	Halobacteriales
15	756883 <i>halophilic archaeon</i> DL31	Euryarchaeota	Halobacteriales
16	797210 <i>Halopiger xanaduensis</i> SH-6	Euryarchaeota	Halobacteriales
17	362976 <i>Haloquadratum walsbyi</i> DSM 16790	Euryarchaeota	Halobacteriales
18	1033806 <i>Halorhabdus tiamatea</i> SARL4B	Euryarchaeota	Halobacteriales
19	519442 <i>Halorhabdus utahensis</i> DSM 12940	* Euryarchaeota	Halobacteriales
20	416348 <i>Halorubrum lacusprofundi</i> ATCC 49239	* Euryarchaeota	Halobacteriales
21	543526 <i>Haloterrigena turkmenica</i> DSM 5511	Euryarchaeota	Halobacteriales
22	547559 <i>Natrialba magadii</i> ATCC 43099	* Euryarchaeota	Halobacteriales
23	797303 <i>Natrinema pellirubrum</i> DSM 15624	Euryarchaeota	Halobacteriales
24	797304 <i>Natronobacterium gregoryi</i> SP2	Euryarchaeota	Halobacteriales
25	348780 <i>Natronomonas pharaonis</i> DSM 2160	* Euryarchaeota	Halobacteriales
26	1072681 <i>Candidatus</i> Haloredivivus sp. G17	* Euryarchaeota	Nanohaloarchaea
27	889948 <i>Candidatus</i> Nanosalina sp. J07AB43	* Euryarchaeota	Nanohaloarchaea
28	889962 <i>Candidatus</i> Nanosalinarum sp. J07AB56	* Euryarchaeota	Nanohaloarchaea
29	410358 <i>Methanocorpusculum labreanum</i> Z	* Euryarchaeota	Methanomicrobiales
30	368407 <i>Methanoculleus marisnigri</i> JR1	* Euryarchaeota	Methanomicrobiales
31	882090 <i>Methanolinea tarda</i> NOBI-1	Euryarchaeota	Methanomicrobiales
32	937775 <i>Methanoplanus limicola</i> DSM 2279	* Euryarchaeota	Methanomicrobiales
33	679926 <i>Methanoplanus petrolearius</i> DSM 11571	Euryarchaeota	Methanomicrobiales
34	456442 <i>Methanoregula boonei</i> 6A8	Euryarchaeota	Methanomicrobiales
35	521011 <i>Methanosphaerula palustris</i> E1-9c	* Euryarchaeota	Methanomicrobiales
36	323259 <i>Methanospirillum hungatei</i> JF-1	* Euryarchaeota	Methanomicrobiales
37	259564 <i>Methanococcoides burtonii</i> DSM 6242	* Euryarchaeota	Methanosarcinales
38	644295 <i>Methanohalobium evestigatum</i> Z-7303	Euryarchaeota	Methanosarcinales
39	547558 <i>Methanohalophilus mahii</i> DSM 5219	Euryarchaeota	Methanosarcinales
40	990316 <i>Methanosaeta concilii</i> GP6	Euryarchaeota	Methanosarcinales
41	1110509 <i>Methanosaeta harundinacea</i> 6Ac	* Euryarchaeota	Methanosarcinales
42	349307 <i>Methanosaeta thermophila</i> PT	* Euryarchaeota	Methanosarcinales
43	679901 <i>Methanosalum zhilinae</i> DSM 4017	* Euryarchaeota	Methanosarcinales
44	188937 <i>Methanosarcina acetivorans</i> C2A	Euryarchaeota	Methanosarcinales
45	269797 <i>Methanosarcina barkeri</i> str. <i>Fusaro</i>	Euryarchaeota	Methanosarcinales
46	192952 <i>Methanosarcina mazei</i> Go1	* Euryarchaeota	Methanosarcinales
47	351160 <i>Methanocella arvoryzae</i> MRE50	* Euryarchaeota	Methanocellales
48	1041930 <i>Methanocella conradii</i> HZ254	* Euryarchaeota	Methanocellales
49	304371 <i>Methanocella paludicola</i> SANAE	* Euryarchaeota	Methanocellales
50	224325 <i>Archaeoglobus fulgidus</i> DSM 4304	* Euryarchaeota	Archaeoglobales
51	572546 <i>Archaeoglobus profundus</i> DSM 5631	* Euryarchaeota	Archaeoglobales
52	693661 <i>Archaeoglobus veneficus</i> SNP6	* Euryarchaeota	Archaeoglobales
53	589924 <i>Ferroglobus placidus</i> DSM 10642	* Euryarchaeota	Archaeoglobales
54	439481 <i>Aciduliprofundum boonei</i> T469	* Euryarchaeota	Thermoplasmatales
55	333146 <i>Ferroplasma acidarmanus</i> fer1	* Euryarchaeota	Thermoplasmatales
56	263820 <i>Picrophilus torridus</i> DSM 9790	* Euryarchaeota	Thermoplasmatales
57	273075 <i>Thermoplasma acidophilum</i> DSM 1728	* Euryarchaeota	Thermoplasmatales
58	273116 <i>Thermoplasma volcanium</i> GSS1	* Euryarchaeota	Thermoplasmatales
59	573064 <i>Methanocaldococcus fervens</i> AG86	Euryarchaeota	Methanococcales
60	573063 <i>Methanocaldococcus infernus</i> ME	* Euryarchaeota	Methanococcales
61	243232 <i>Methanocaldococcus jannaschii</i> DSM 2661	* Euryarchaeota	Methanococcales
62	644281 <i>Methanocaldococcus</i> sp. FS406-22	Euryarchaeota	Methanococcales
63	579137 <i>Methanocaldococcus vulcanius</i> M7	Euryarchaeota	Methanococcales
64	419665 <i>Methanococcus aeolicus</i> Nankai-3	* Euryarchaeota	Methanococcales
65	402880 <i>Methanococcus maripaludis</i> C5	Euryarchaeota	Methanococcales
66	406327 <i>Methanococcus vannieli</i> SB	* Euryarchaeota	Methanococcales

67	456320	<i>Methanococcus voltae</i> A3	Euryarchaeota	Methanococcales
68	647113	<i>Methanothermococcus okinawensis</i> IH1	Euryarchaeota	Methanococcales
69	647171	<i>Methanotorris formicicus</i> Mc-S-70	Euryarchaeota	Methanococcales
70	880724	<i>Methanotorris igneus</i> Kol 5	* Euryarchaeota	Methanococcales
71	868132	<i>Methanobacterium</i> sp. AL-21	* Euryarchaeota	Methanobacteriales
72	634498	<i>Methanobrevibacter ruminantium</i> M1	Euryarchaeota	Methanobacteriales
73	420247	<i>Methanobrevibacter smithii</i> ATCC 35061	* Euryarchaeota	Methanobacteriales
74	339860	<i>Methanosphaera stadthanae</i> DSM 3091	* Euryarchaeota	Methanobacteriales
75	79929	<i>Methanothermobacter marburgensis</i> str. Marburg	Euryarchaeota	Methanobacteriales
76	187420	<i>Methanothermobacter thermautotrophicus</i> str. Delta h	* Euryarchaeota	Methanobacteriales
77	523846	<i>Methanothermus fervidus</i> DSM 2088	* Euryarchaeota	Methanobacteriales
78	190192	<i>Methanopyrus kandleri</i> AV19	* Euryarchaeota	Methanopyrales
79	272844	<i>Pyrococcus abyssii</i> GE5	* Euryarchaeota	Thermococcales
80	186497	<i>Pyrococcus furiosus</i> DSM 3638	Euryarchaeota	Thermococcales
81	70601	<i>Pyrococcus horikoshii</i> OT3	Euryarchaeota	Thermococcales
82	342949	<i>Pyrococcus</i> sp. NA2	Euryarchaeota	Thermococcales
83	529709	<i>Pyrococcus yayanosii</i> CH1	* Euryarchaeota	Thermococcales
84	391623	<i>Thermococcus barophilus</i> MP	* Euryarchaeota	Thermococcales
85	593117	<i>Thermococcus gammatolerans</i> EJ3	* Euryarchaeota	Thermococcales
86	69014	<i>Thermococcus kodakarensis</i> KOD1	Euryarchaeota	Thermococcales
87	523849	<i>Thermococcus litoralis</i> DSM 5473	* Euryarchaeota	Thermococcales
88	523850	<i>Thermococcus onnurineus</i> NA1	Euryarchaeota	Thermococcales
89	604354	<i>Thermococcus sibiricus</i> MM 739	Euryarchaeota	Thermococcales
90	246969	<i>Thermococcus</i> sp. AM4	Euryarchaeota	Thermococcales
91	425595	<i>Candidatus</i> Micrarchaeum acidiphilum ARMAN-2	* Euryarchaeota	ARMAN
92	662760	<i>Candidatus</i> Parvarchaeum acidiphilum ARMAN-4	* Euryarchaeota	ARMAN
93	662762	<i>Candidatus</i> Parvarchaeum acidophilus ARMAN-5	* Euryarchaeota	ARMAN
94	228908	<i>Nanoarchaeum equitans</i> Kin4-M	* Euryarchaeota	Nanoarchaeota
95	274854	uncultured marine group II euryarchaeote	* Euryarchaeota	Unclassified
96	666510	<i>Acidilobus saccharovorans</i> 345-15	* Crenarchaeota	Desulfurococcales
97	272557	<i>Aeropyrum pernix</i> K1	* Crenarchaeota	Desulfurococcales
98	768672	<i>Desulfurococcus fermentans</i> DSM 16532	Crenarchaeota	Desulfurococcales
99	490899	<i>Desulfurococcus kamchatkensis</i> 1221n	* Crenarchaeota	Desulfurococcales
100	765177	<i>Desulfurococcus mucosus</i> DSM 2162	Crenarchaeota	Desulfurococcales
101	415426	<i>Hyperthermus butylicus</i> DSM 5456	* Crenarchaeota	Desulfurococcales
102	453591	<i>Ignicoccus hospitalis</i> KIN4/I	* Crenarchaeota	Desulfurococcales
103	583356	<i>Ignisphaera aggregans</i> DSM 17230	* Crenarchaeota	Desulfurococcales
104	694429	<i>Pyrolobus fumarii</i> 1A	* Crenarchaeota	Desulfurococcales
105	591019	<i>Staphylothermus hellenicus</i> DSM 12710	* Crenarchaeota	Desulfurococcales
106	399550	<i>Staphylothermus marinus</i> F1	Crenarchaeota	Desulfurococcales
107	633148	<i>Thermosphaera aggregans</i> DSM 11486	* Crenarchaeota	Desulfurococcales
108	933801	<i>Acidianus hospitalis</i> W1	* Crenarchaeota	Sulfolobales
109	1006006	<i>Metallosphaera cuprina</i> Ar-4	Crenarchaeota	Sulfolobales
110	399549	<i>Metallosphaera sedula</i> DSM 5348	* Crenarchaeota	Sulfolobales
111	671065	<i>Metallosphaera yellowstonensis</i> MK1	* Crenarchaeota	Sulfolobales
112	330779	<i>Sulfolobus acidocaldarius</i> DSM 639	* Crenarchaeota	Sulfolobales
113	429572	<i>Sulfolobus islandicus</i> L.S.2.15	Crenarchaeota	Sulfolobales
114	273057	<i>Sulfolobus solfataricus</i> P2	* Crenarchaeota	Sulfolobales
115	273063	<i>Sulfolobus tokodaii</i> str. 7	* Crenarchaeota	Sulfolobales
116	397948	<i>Caldivirga maquilingensis</i> IC-167	* Crenarchaeota	Thermoproteales
117	178306	<i>Pyrobaculum aerophilum</i> str. IM2	* Crenarchaeota	Thermoproteales
118	340102	<i>Pyrobaculum arsenaticum</i> DSM 13514	* Crenarchaeota	Thermoproteales
119	410359	<i>Pyrobaculum caldifontis</i> JCM 11548	Crenarchaeota	Thermoproteales
120	384616	<i>Pyrobaculum islandicum</i> DSM 4184	* Crenarchaeota	Thermoproteales
121	698757	<i>Pyrobaculum neutrophilum</i> V24Sta	Crenarchaeota	Thermoproteales
122	1104324	<i>Pyrobaculum oguniense</i> TE7	Crenarchaeota	Thermoproteales
123	368408	<i>Pyrobaculum</i> sp. 1860	Crenarchaeota	Thermoproteales
124	444157	<i>Thermofilum pendens</i> Hrk 5	* Crenarchaeota	Thermoproteales
125	768679	<i>Thermoproteus tenax</i> Kra 1	* Crenarchaeota	Thermoproteales
126	999630	<i>Thermoproteus uzoniensis</i> 768-20	Crenarchaeota	Thermoproteales
127	572478	<i>Vulcanisaeta distributa</i> DSM 14429	* Crenarchaeota	Thermoproteales
128	985053	<i>Vulcanisaeta moutnovskia</i> 768-28	Crenarchaeota	Thermoproteales
129	374847	<i>Candidatus</i> Korarchaeum cryptofilum OPF8	* Korarchaeota	



Petitjean et al.
 Supplementary Figure S1
 200 PhyML-SH LG+G8

'Ca. Corarchaeum cryptofilum OPF8'

'Ca. Caldarchaeum subterraneum'

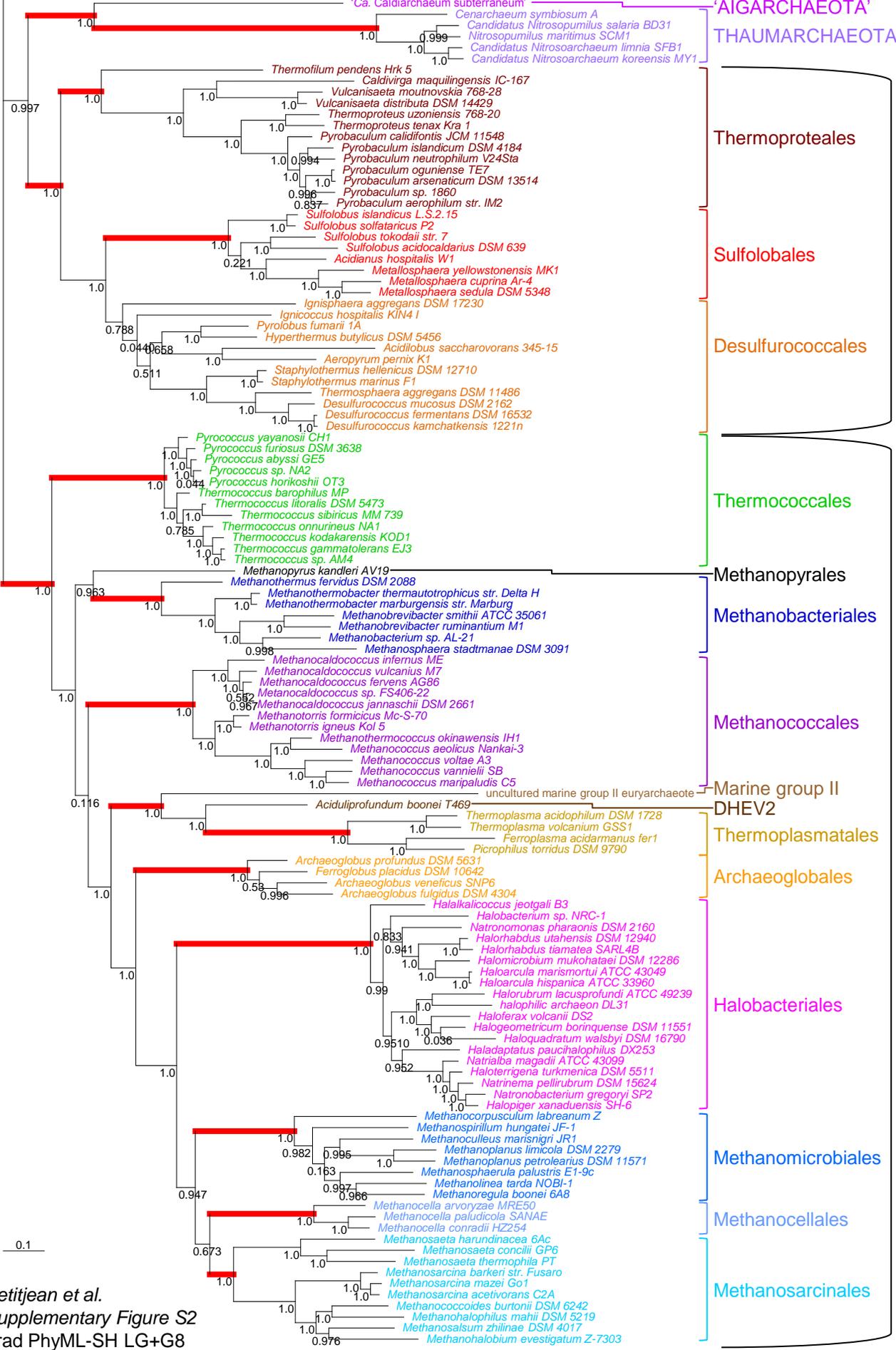
KORARCHAEOTA

'AIGARCHAEOTA'

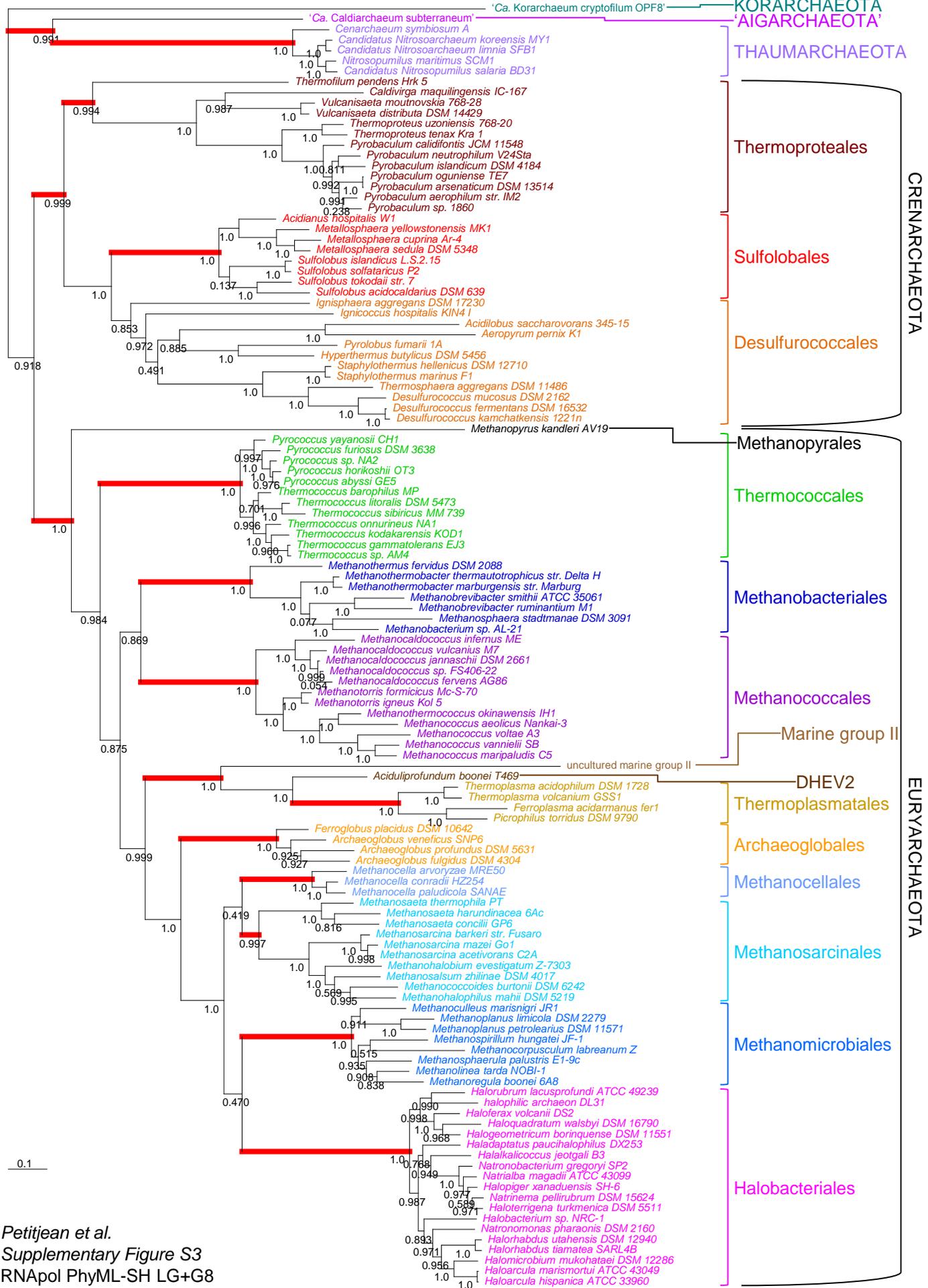
THAUMARCHAEOTA

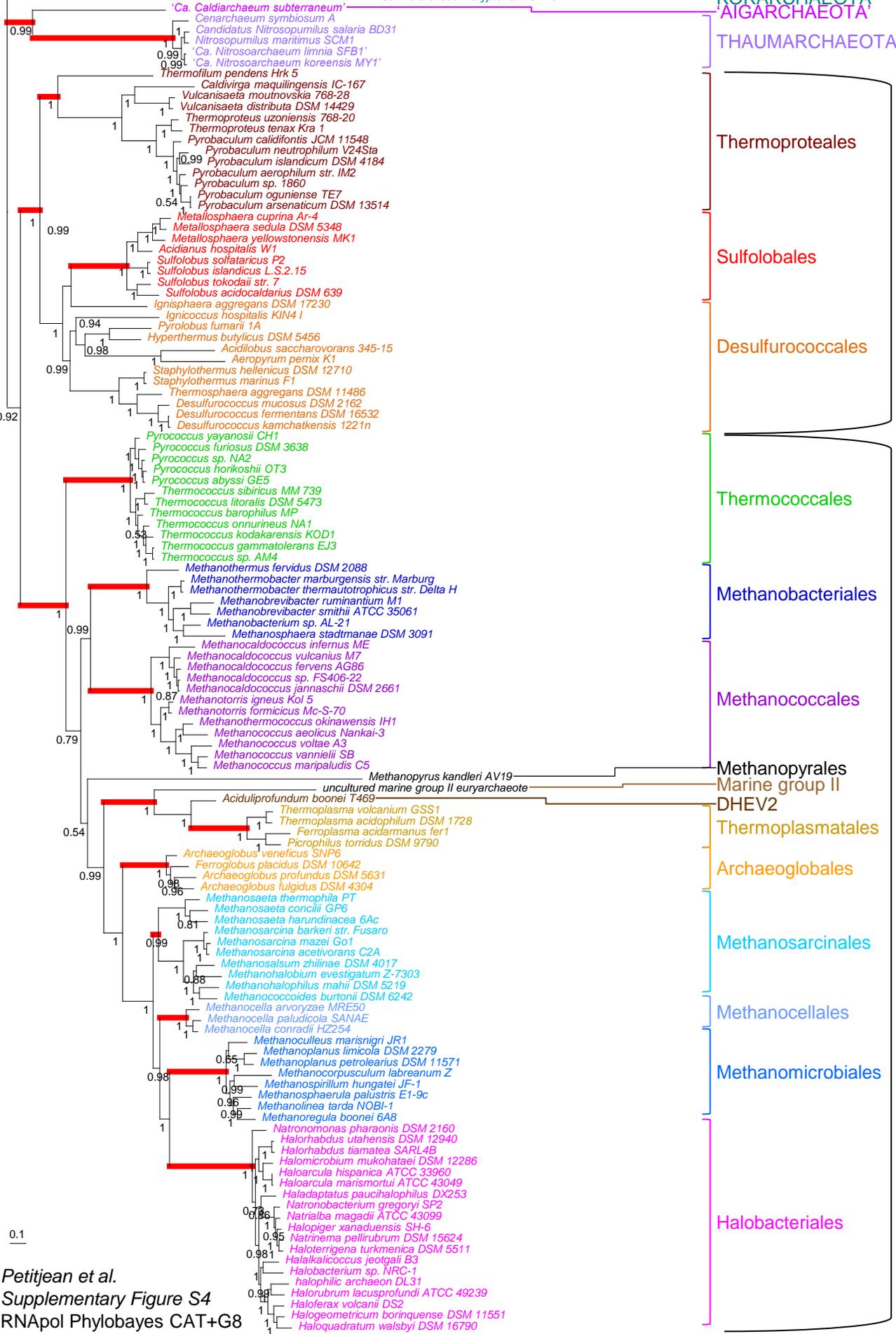
CRENARCHAEOTA

EURYARCHAEOTA

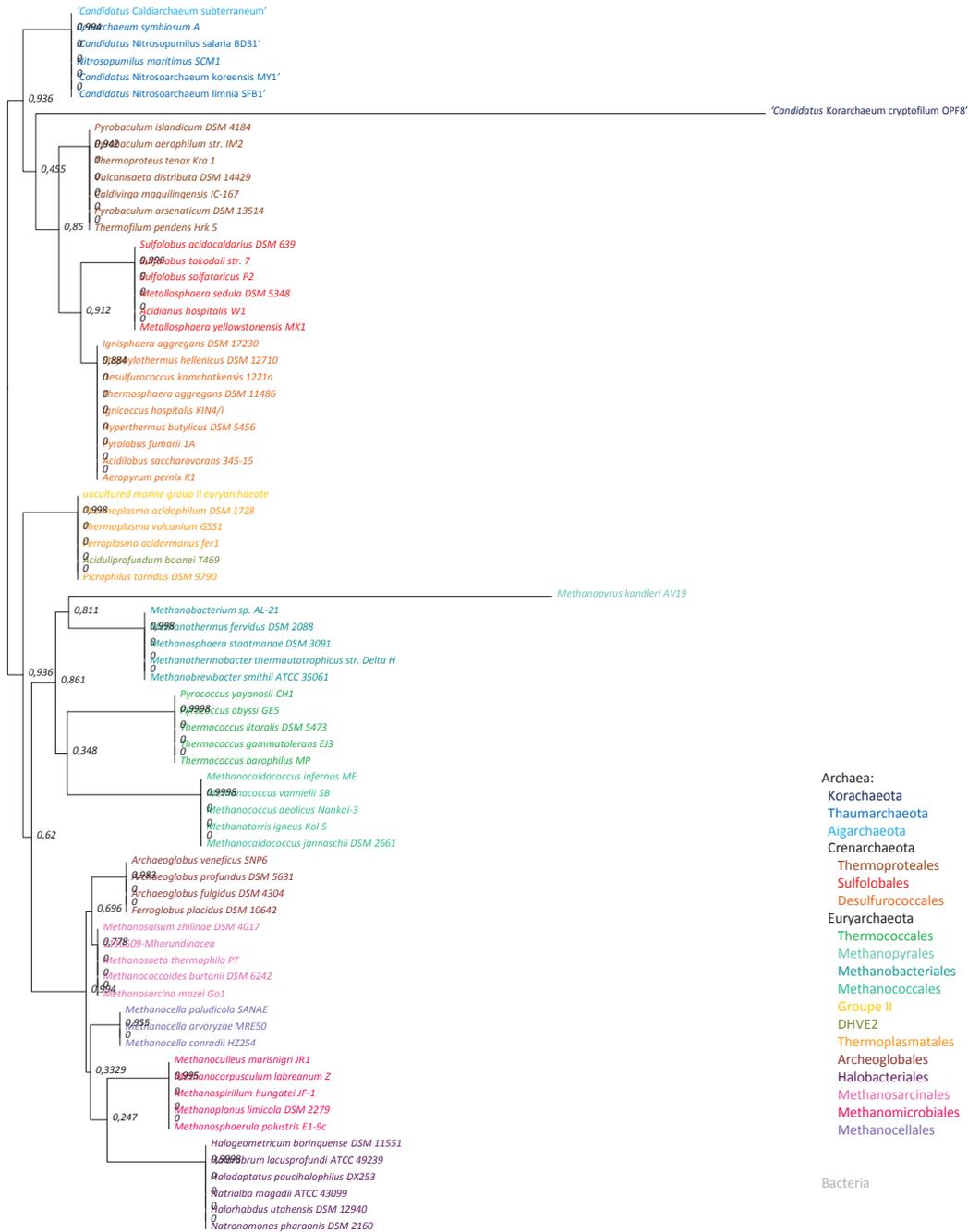


Petitjean et al.
 Supplementary Figure S2
 Trad PhyML-SH LG+G8





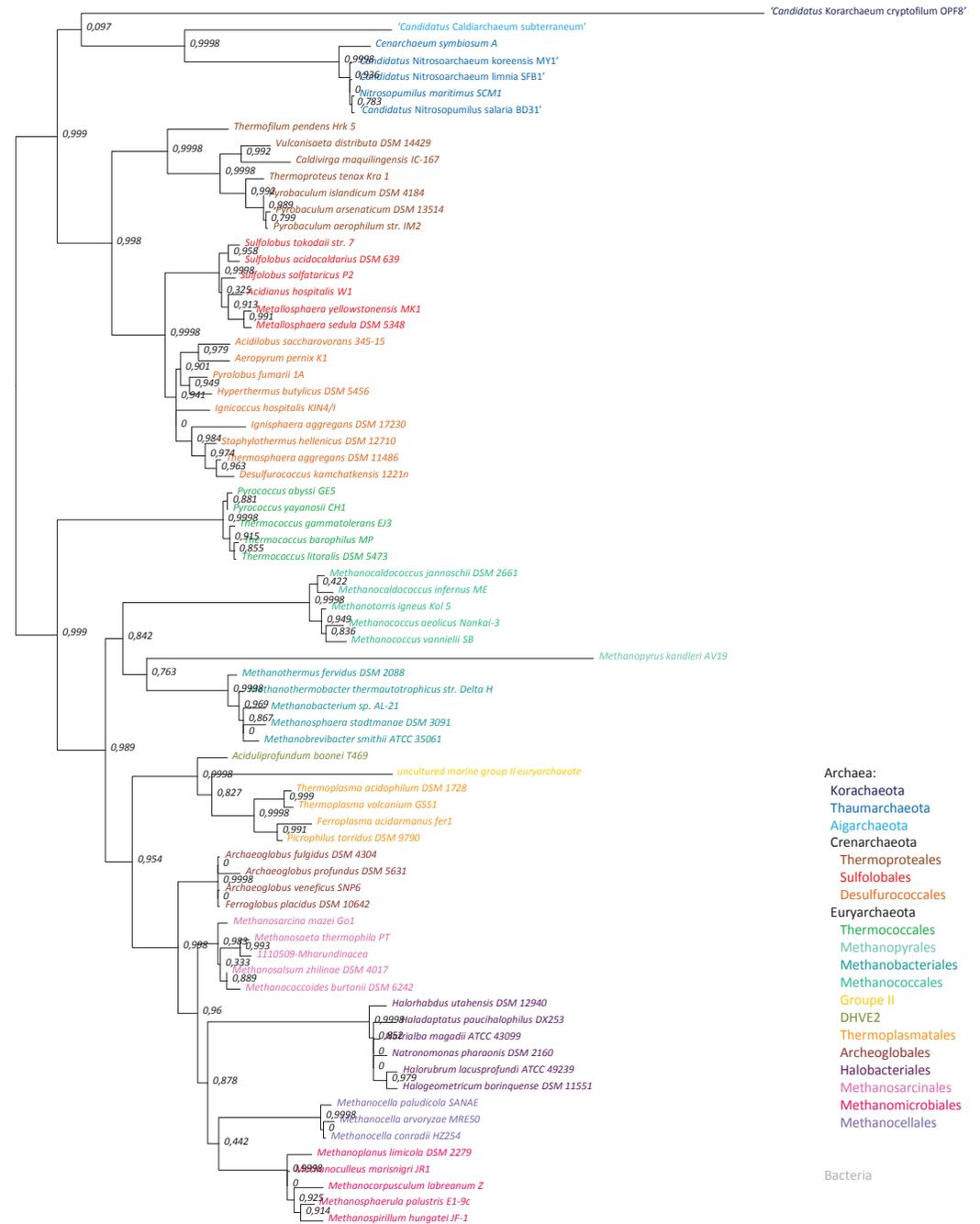
Petitjean et al. Supplementary Figure S4 RNAPol Phylobayes CAT+G8



Petitjean et al.
Supplementary Figure S5
SF₀, PhyML-SH LG+G8

0.01

74 species - 1958 positions



Petitjean et al.
Supplementary Figure S5
SF₁, PhyML-SH LG+G8

0.02

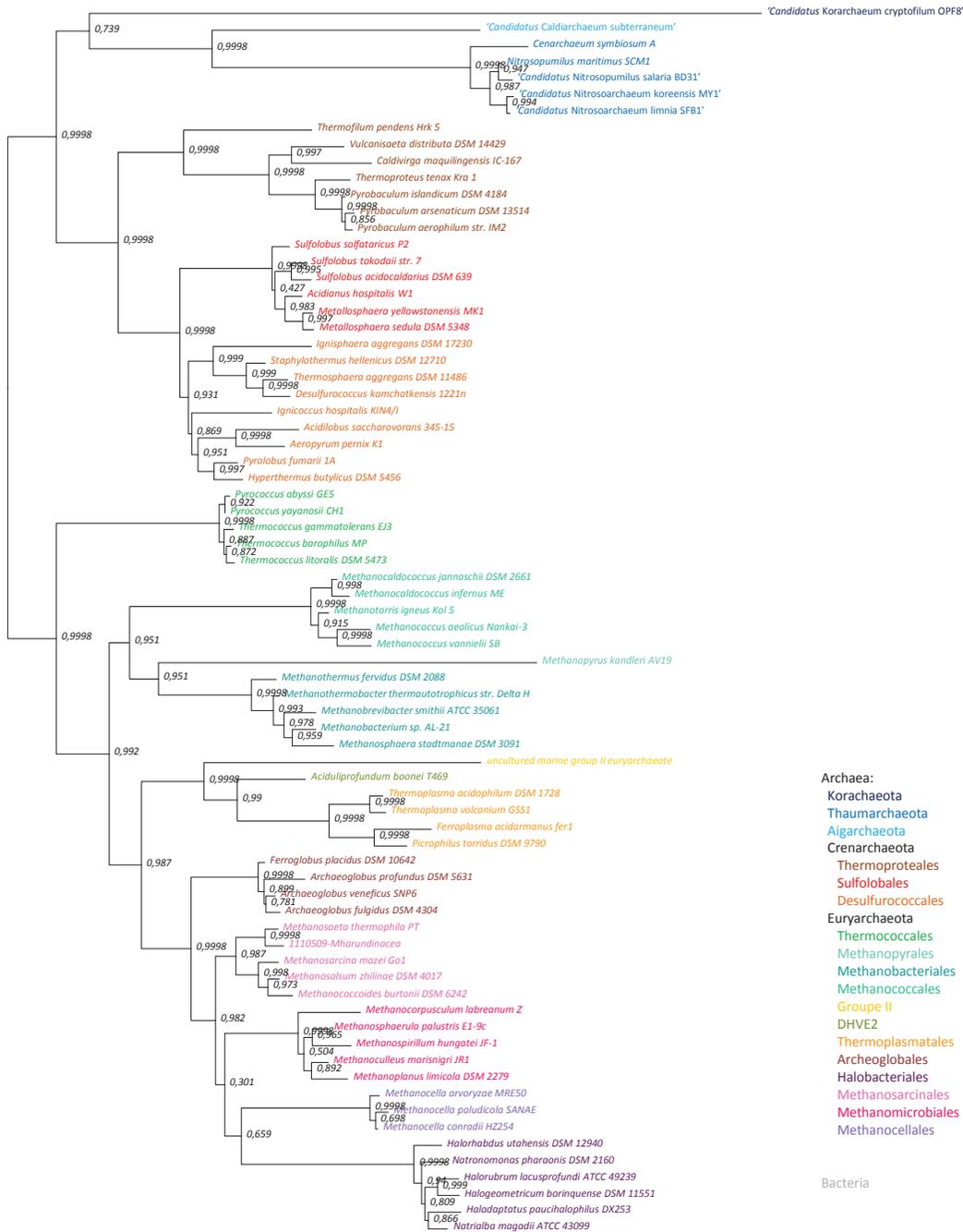
74 species - 3884 positions

Archaea:
Korarchaeota
Thaumarchaeota
Aigarchaeota
Crenarchaeota
Thermoproteales
Sulfobiales
Desulfurococcales
Euryarchaeota
Thermococcales
Methanopyrales
Methanobacterales
Methanococcales
Groupe II
DHVE2
Thermoplasmatales
Archeoglobales
Halobacteriales
Methanosarcinales
Methanomicrobiales
Methanocellales

Bacteria

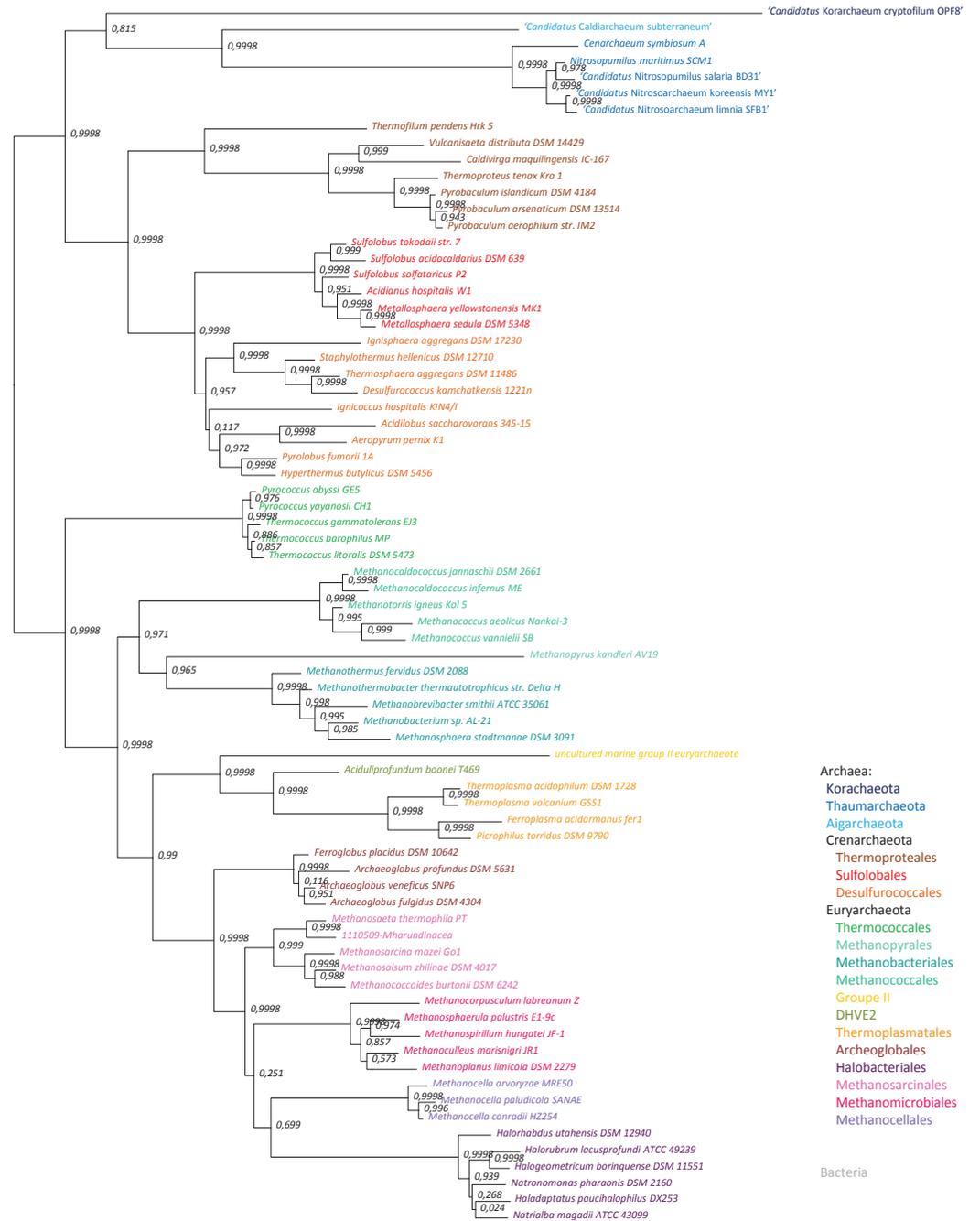
Archaea:
Korarchaeota
Thaumarchaeota
Aigarchaeota
Crenarchaeota
Thermoproteales
Sulfobiales
Desulfurococcales
Euryarchaeota
Thermococcales
Methanopyrales
Methanobacterales
Methanococcales
Groupe II
DHVE2
Thermoplasmatales
Archeoglobales
Halobacteriales
Methanosarcinales
Methanomicrobiales
Methanocellales

Bacteria



Petitjean et al.
 Supplementary Figure S5
 SF₂ PhyML-SH LG+G8

74 species - 5439 positions



Petitjean et al.
 Supplementary Figure S5
 SF₃ PhyML-SH LG+G8

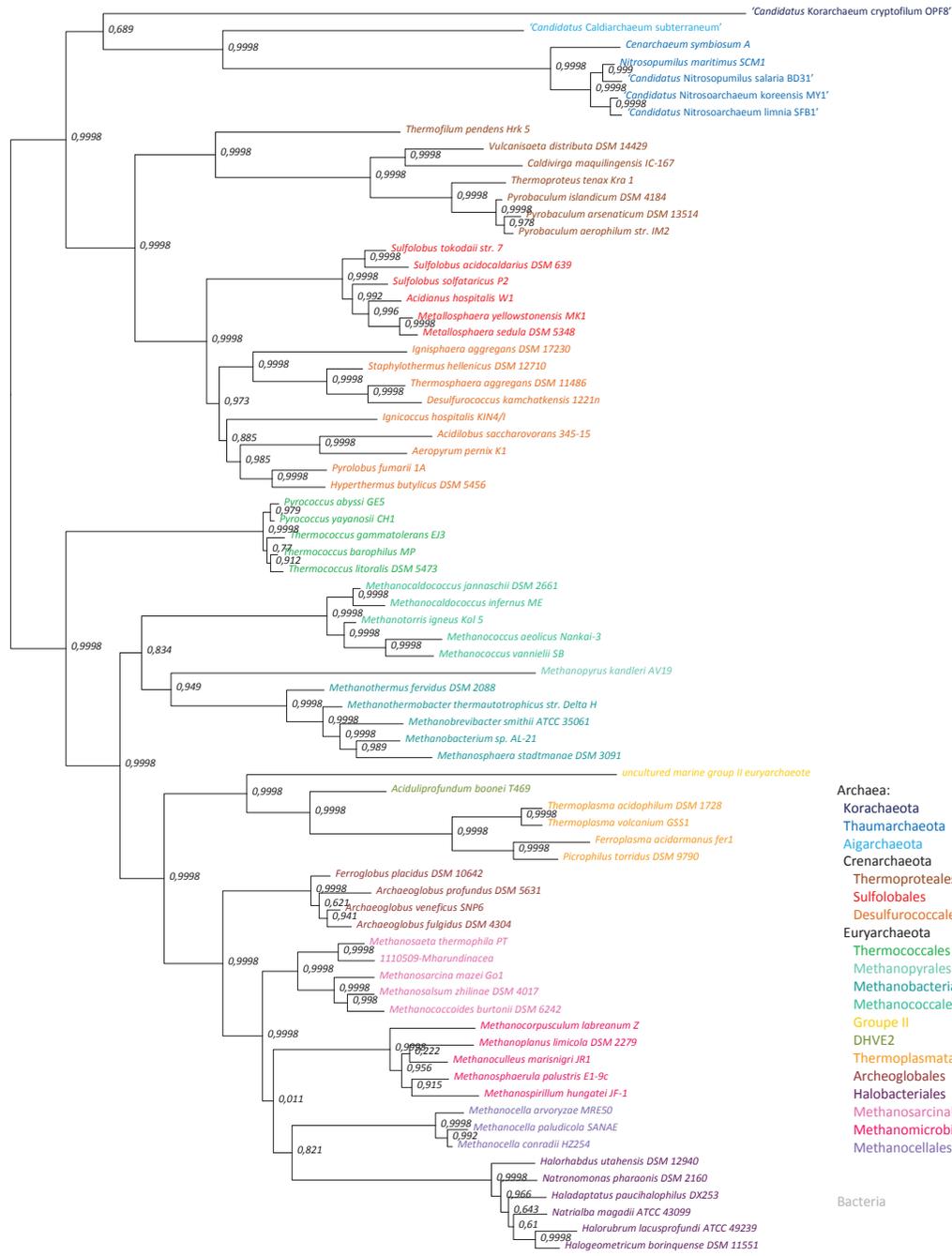
74 species - 6929 positions

Archaea:
 Korarchaeota
 Thaumarchaeota
 Aigarchaeota
 Crenarchaeota
 Thermoproteales
 Sulfolobales
 Desulfurococcales
 Euryarchaeota
 Thermococcales
 Methanopyrales
 Methanobacteriales
 Methanococcales
 Groupe II
 DHVE2
 Thermoplasmatales
 Archeoglobales
 Halobacteriales
 Methanosarcinales
 Methanomicrobiales
 Methanocellales

Bacteria

Archaea:
 Korarchaeota
 Thaumarchaeota
 Aigarchaeota
 Crenarchaeota
 Thermoproteales
 Sulfolobales
 Desulfurococcales
 Euryarchaeota
 Thermococcales
 Methanopyrales
 Methanobacteriales
 Methanococcales
 Groupe II
 DHVE2
 Thermoplasmatales
 Archeoglobales
 Halobacteriales
 Methanosarcinales
 Methanomicrobiales
 Methanocellales

Bacteria



Petitjean et al.
Supplementary Figure S5
SF₄ PhyML-SH LG+G8

0.04

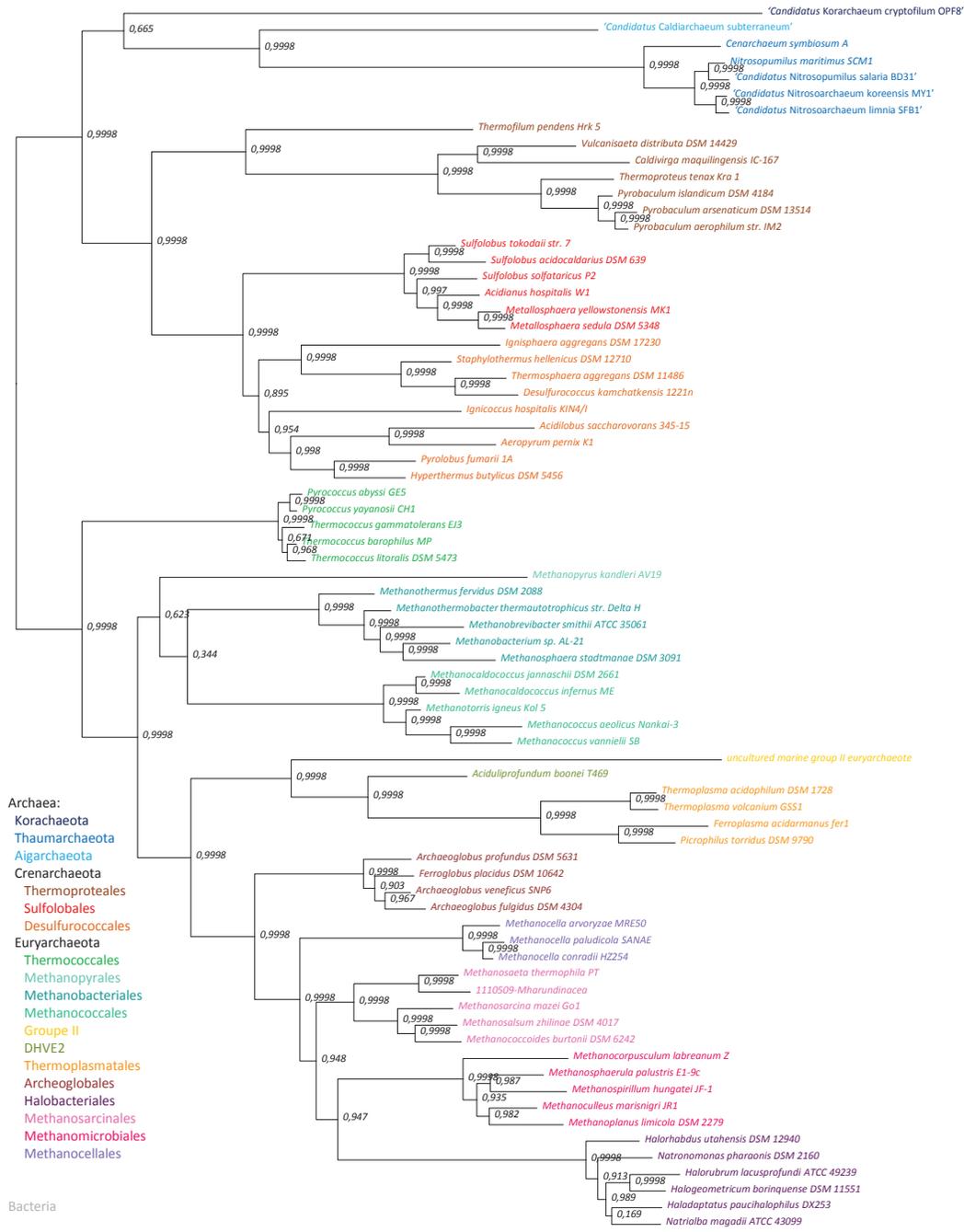
74 species - 8373 positions



Petitjean et al.
Supplementary Figure S5
SF₅ PhyML-SH LG+G8

0.04

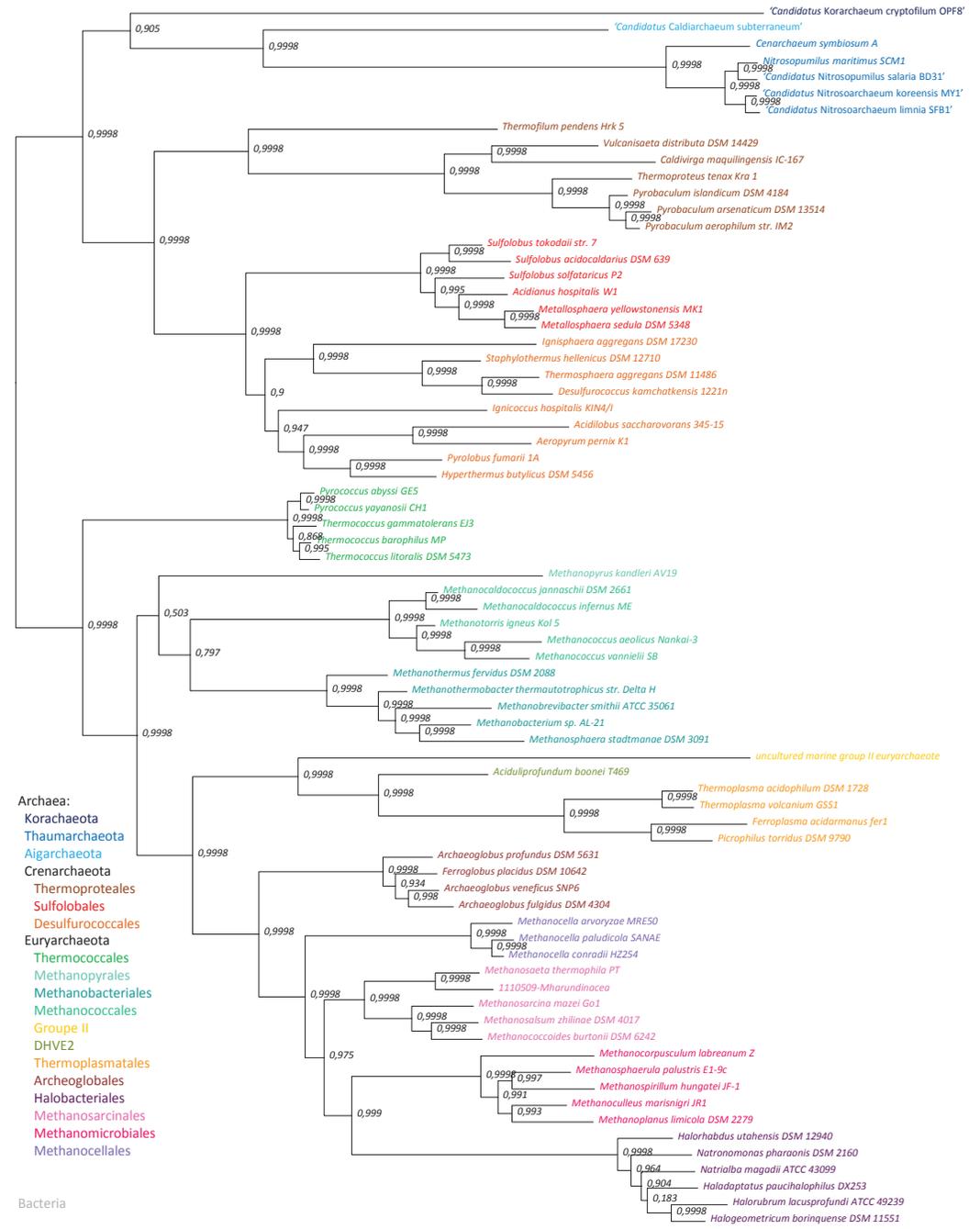
74 species - 9742 positions



Petitjean et al.
Supplementary Figure S5
SF₆ PhyML-SH LG+G8

0.05

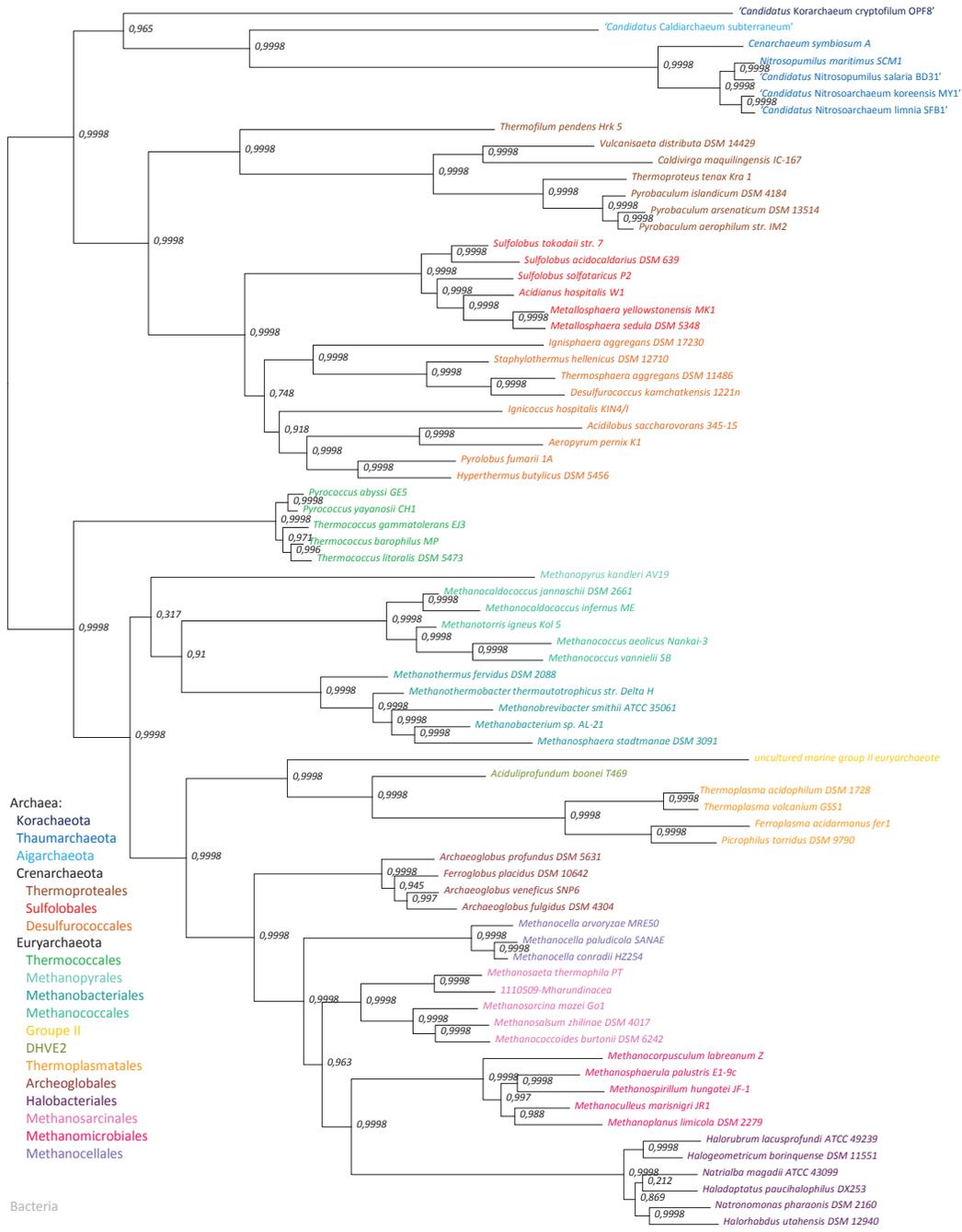
74 species – 11 089 positions



Petitjean et al.
Supplementary Figure S5
SF₇ PhyML-SH LG+G8

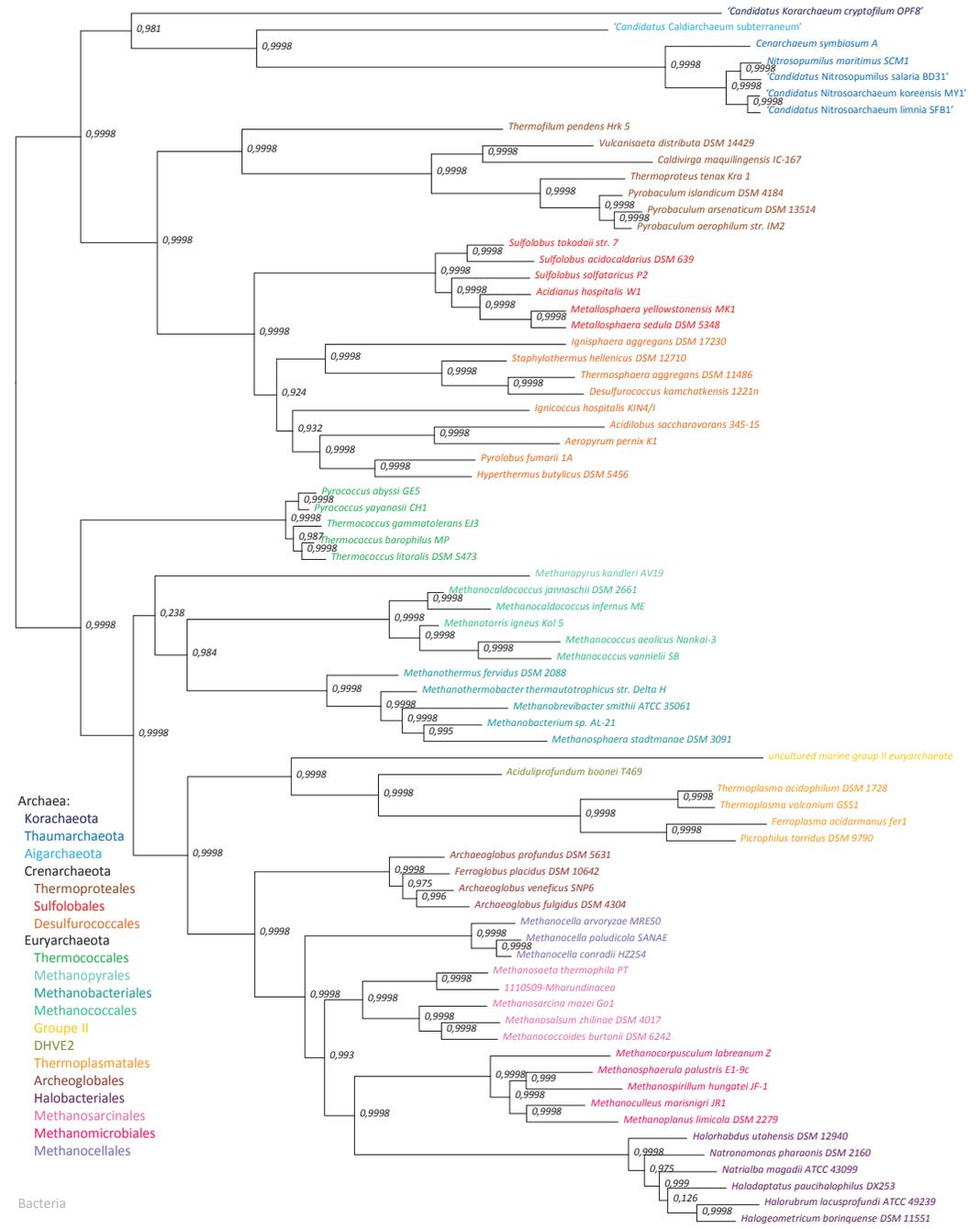
0.06

74 species – 12 465 positions



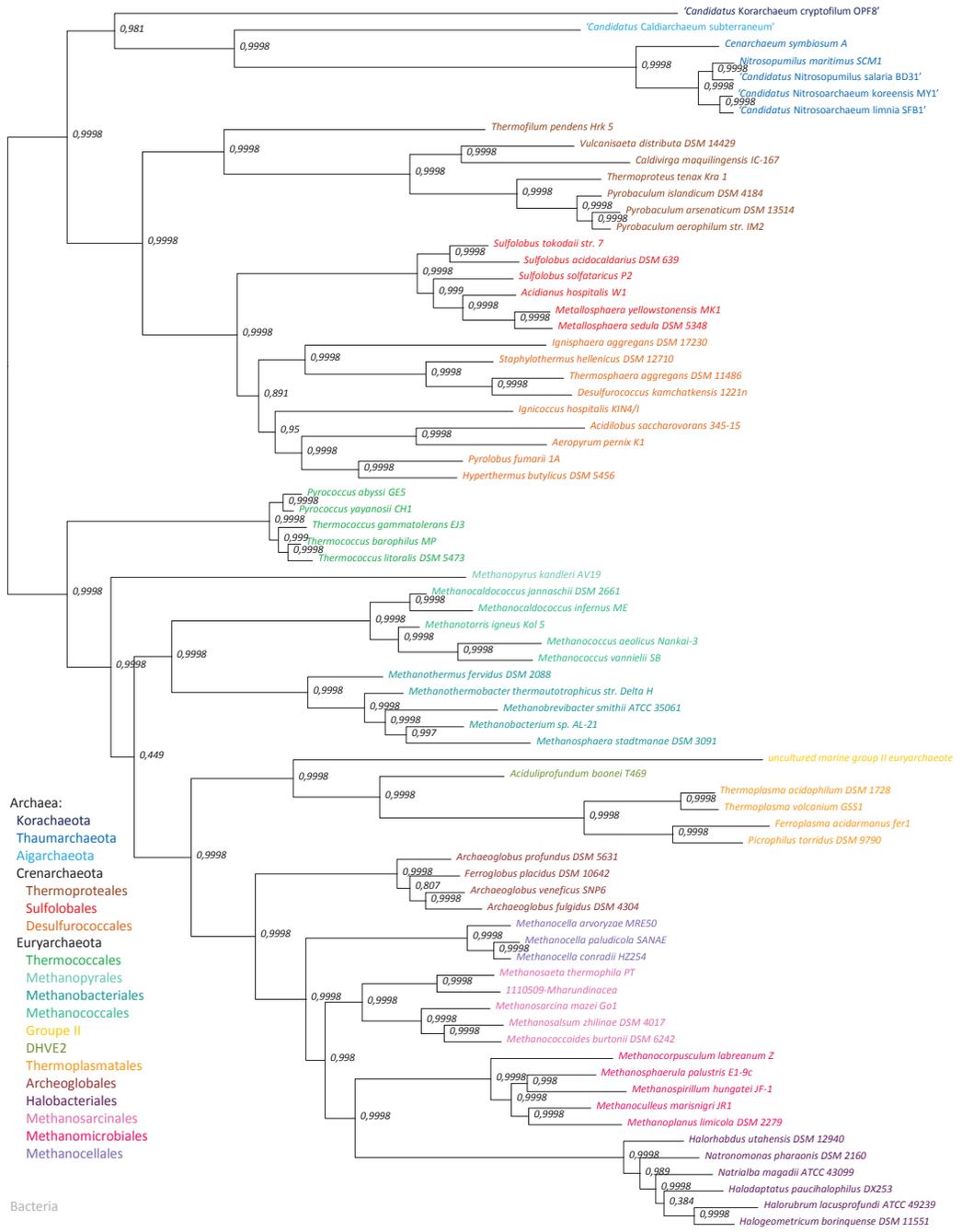
Petitjean et al.
 Supplementary Figure S5
 SF₈ PhyML-SH LG+G8

74 species – 13 806 positions



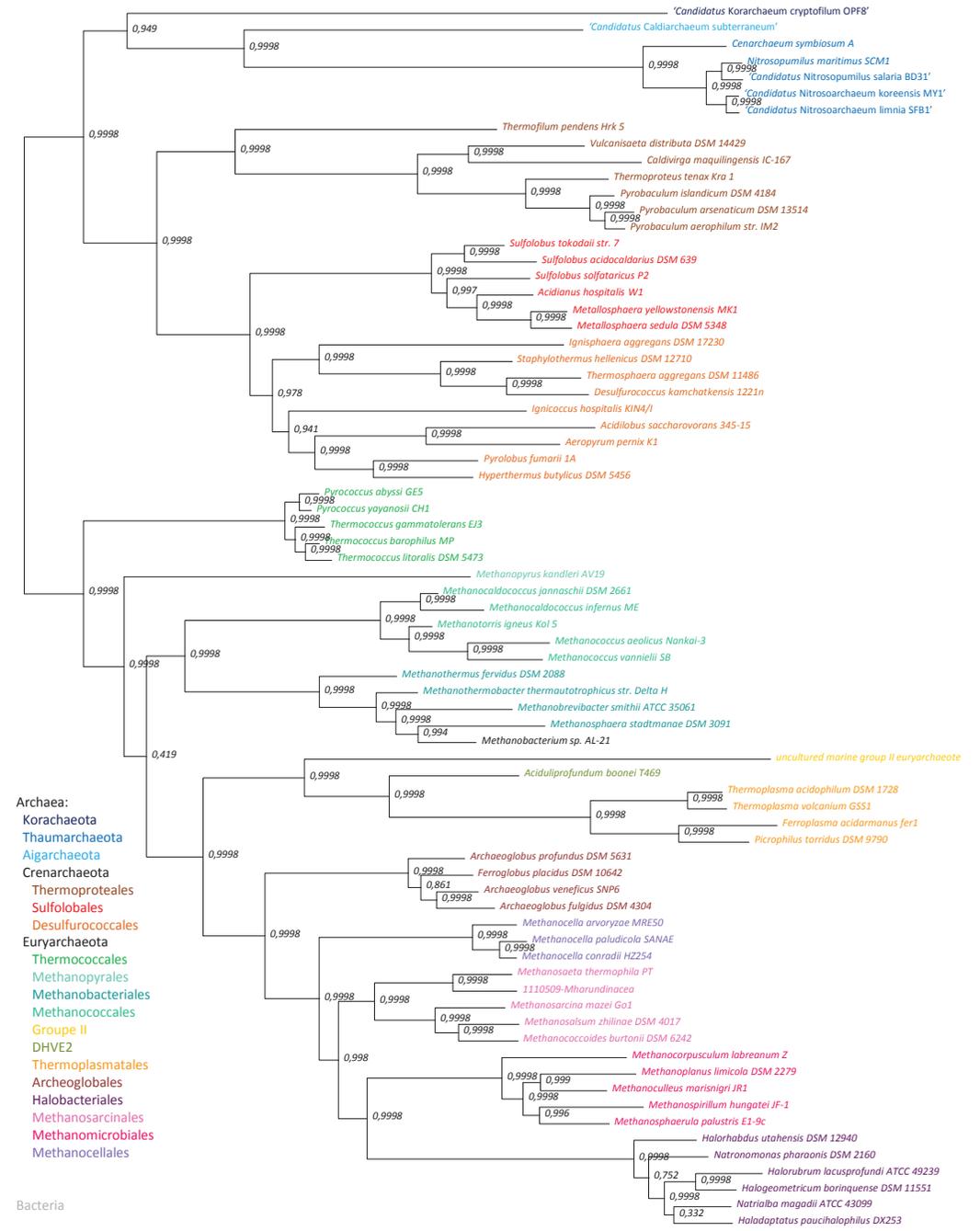
Petitjean et al.
 Supplementary Figure S5
 SF₉ PhyML-SH LG+G8

74 species – 15 176 positions



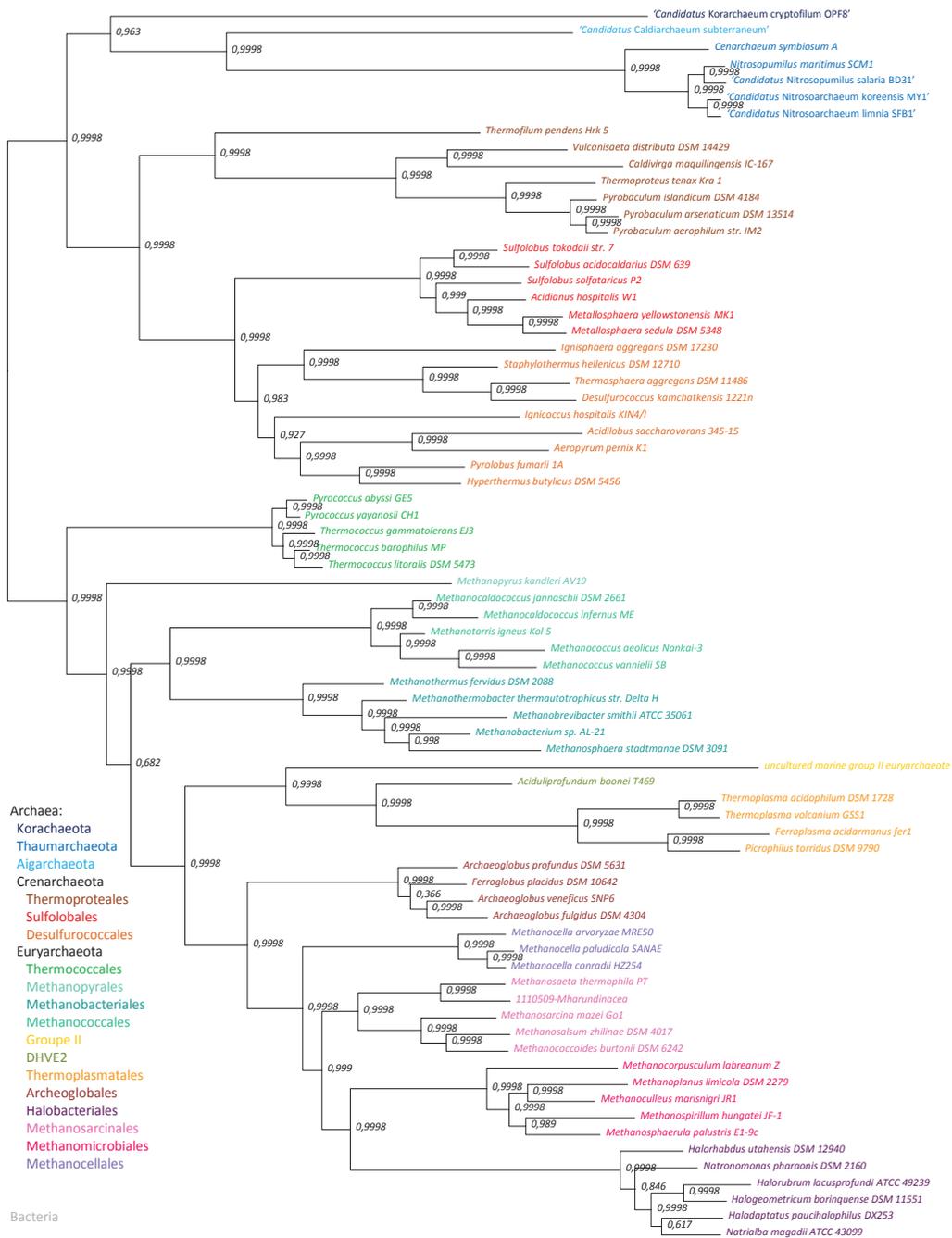
Petitjean et al.
 Supplementary Figure S5
 SF₁₀ PhyML-SH LG+G8

74 species – 16 508 positions



Petitjean et al.
 Supplementary Figure S5
 SF₁₁ PhyML-SH LG+G8

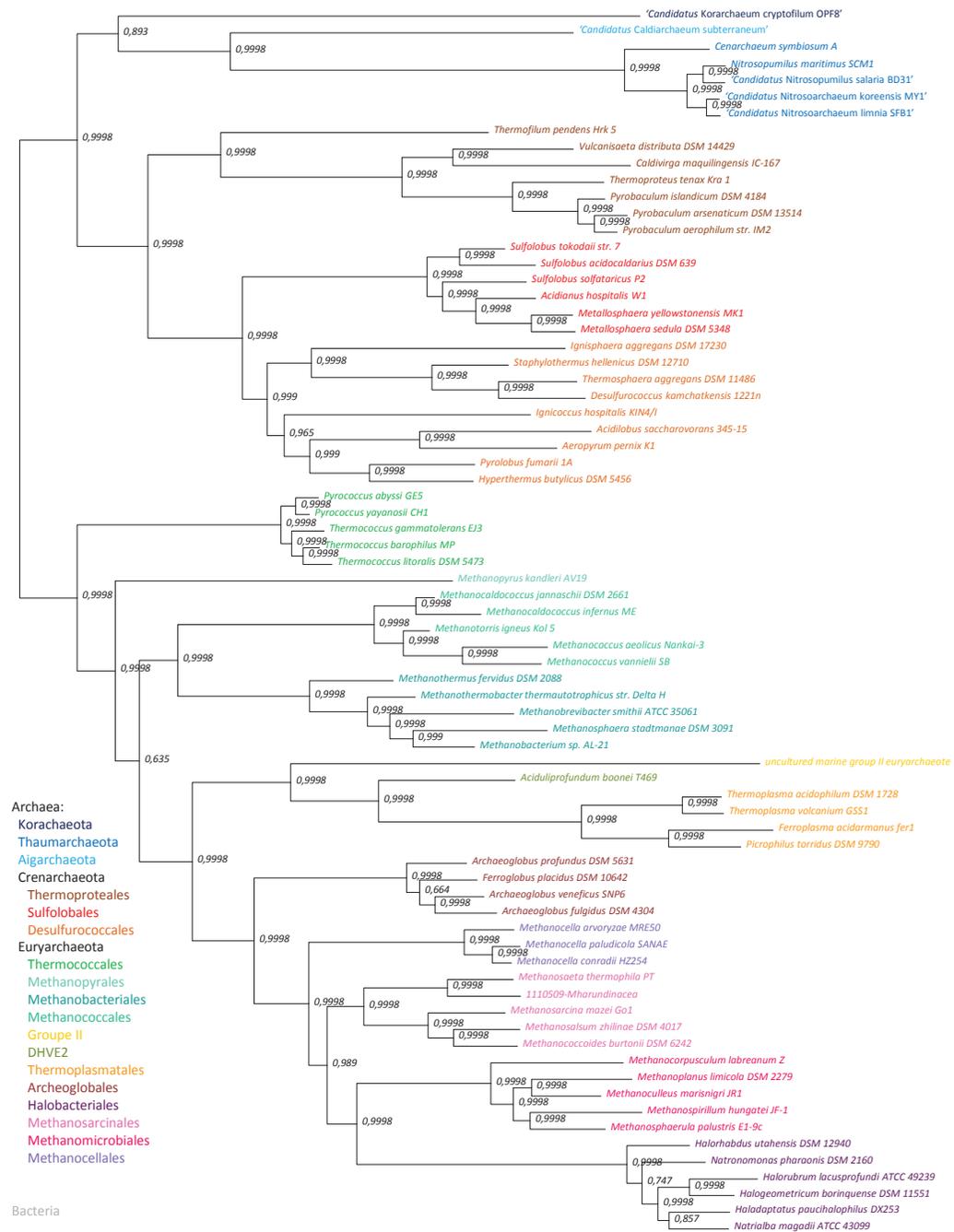
74 species – 17 845 positions



Petitjean et al.
 Supplementary Figure S5
 SF₁₂ PhyML-SH LG+G8

0.1

74 species – 19 167 positions



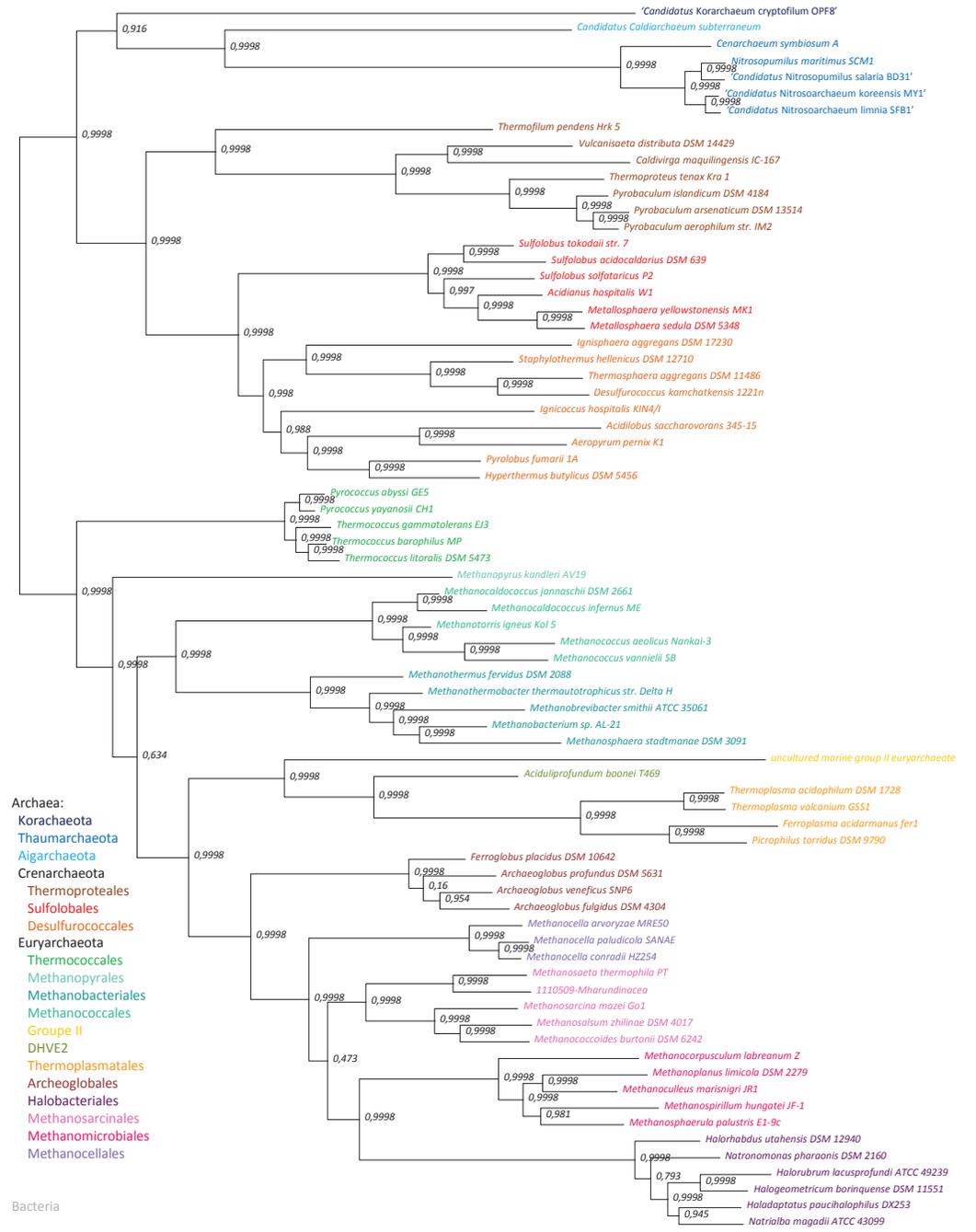
Petitjean et al.
 Supplementary Figure S5
 SF₁₂ PhyML-SH LG+G8

0.2

74 species – 20 458 positions



Petitjean et al.
 Supplementary Figure S5
 SF₁₄ PhyML-SH LG+G8



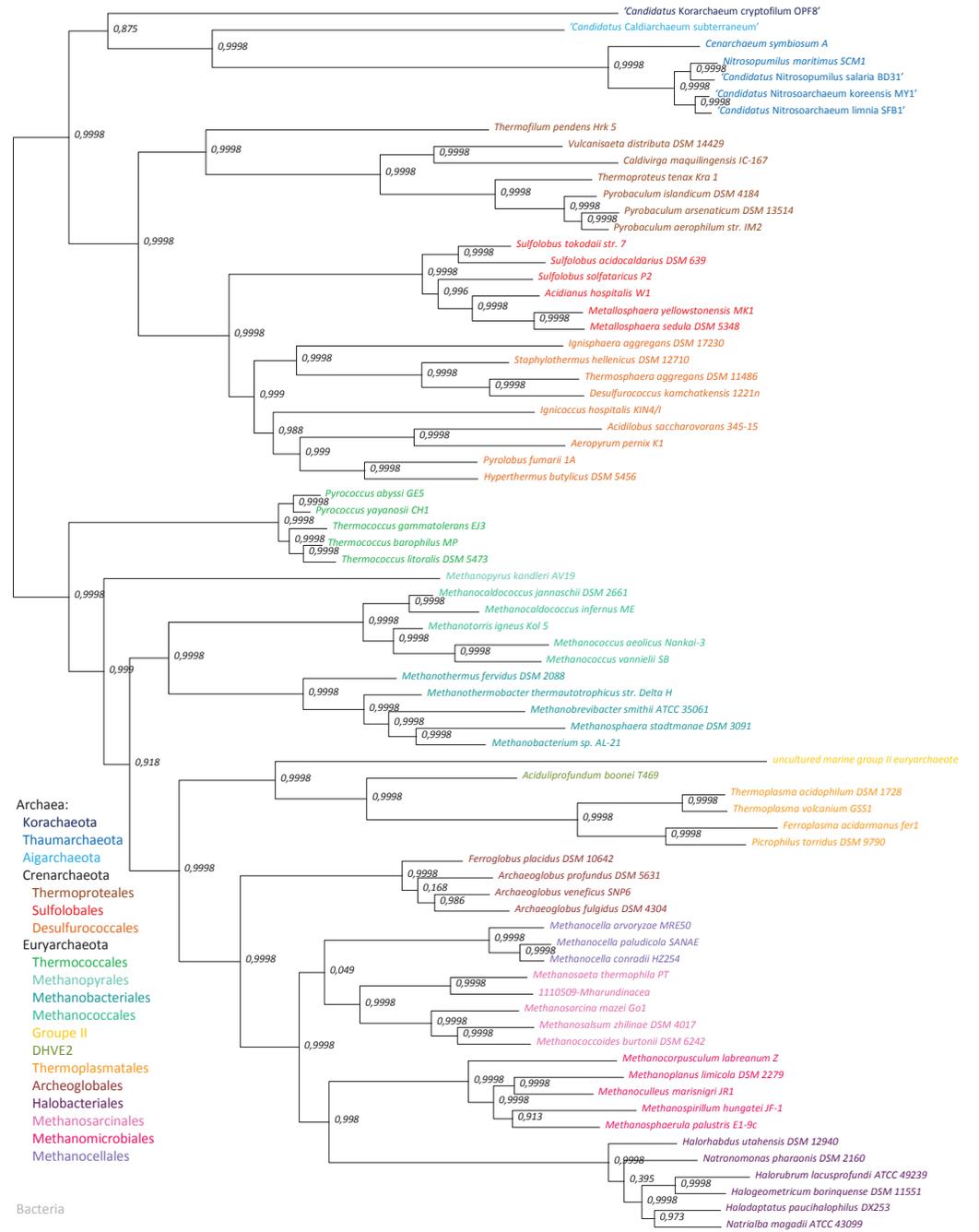
Petitjean et al.
 Supplementary Figure S5
 SF₁₅ PhyML-SH LG+G8



Petitjean et al.
 Supplementary Figure S5
 SF₁₆ PhyML-SH LG+G8

0.2

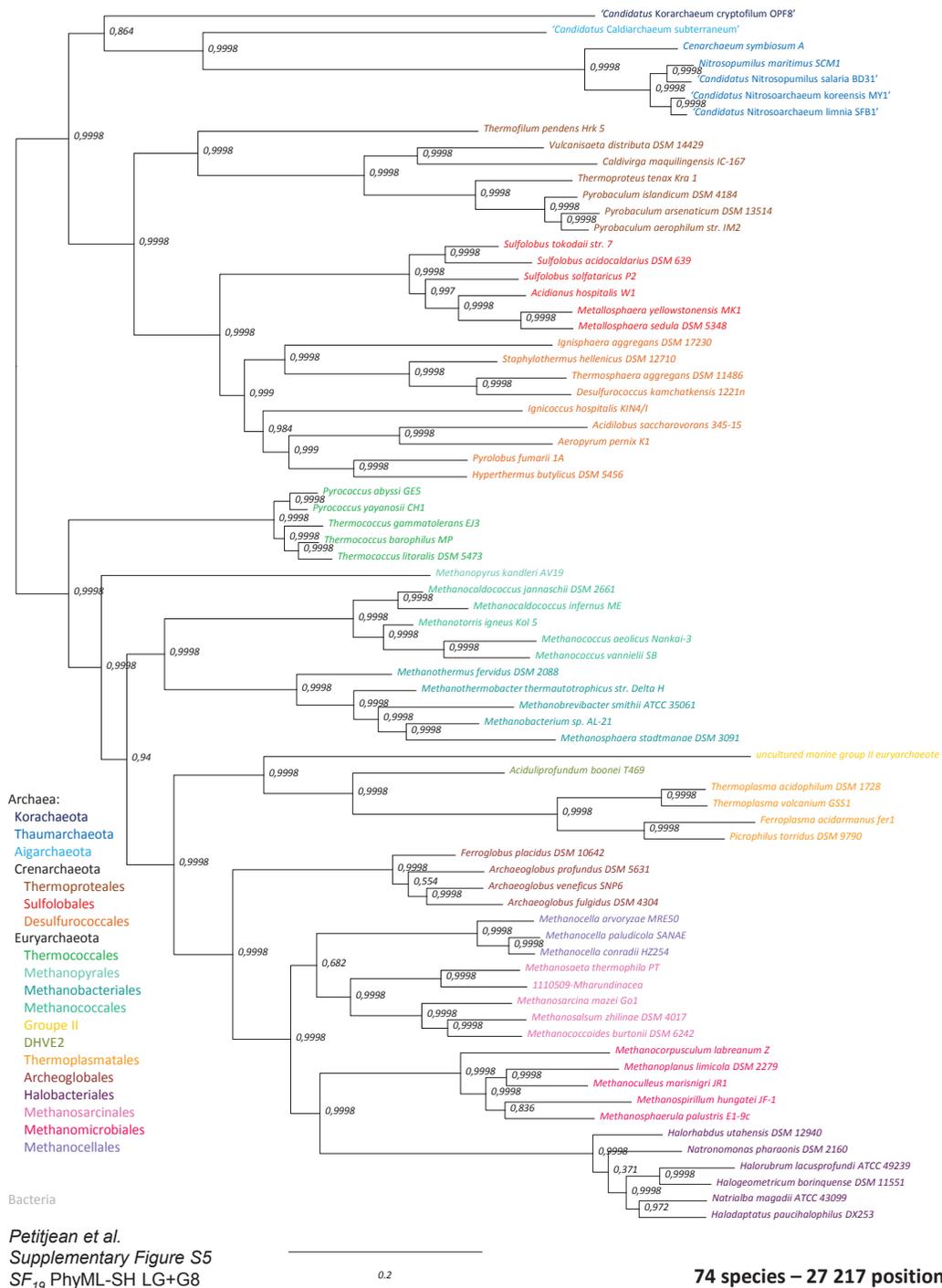
74 species – 24 143 positions



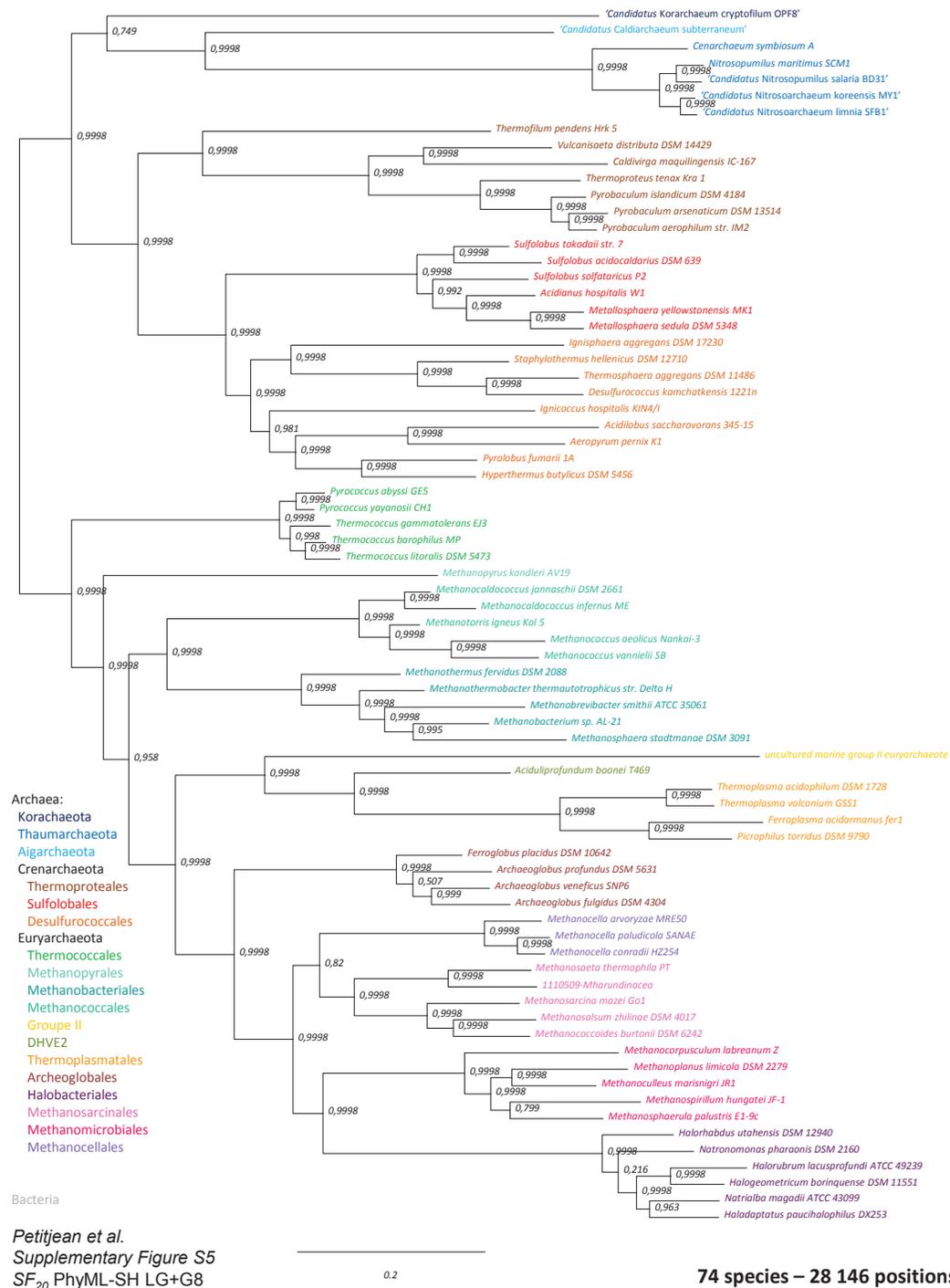
Petitjean et al.
 Supplementary Figure S5
 SF₁₇ PhyML-SH LG+G8

0.2

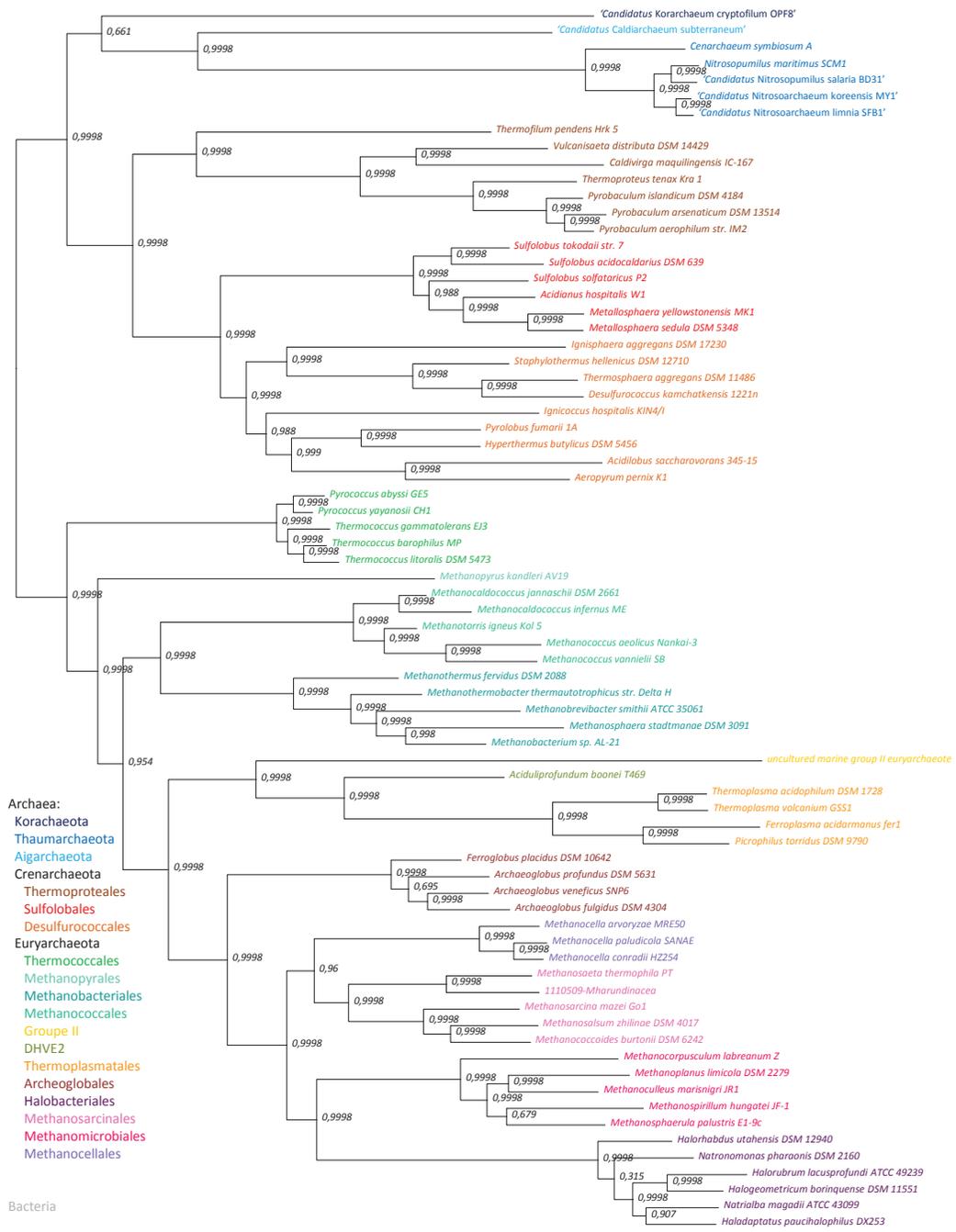
74 species – 25 237 positions



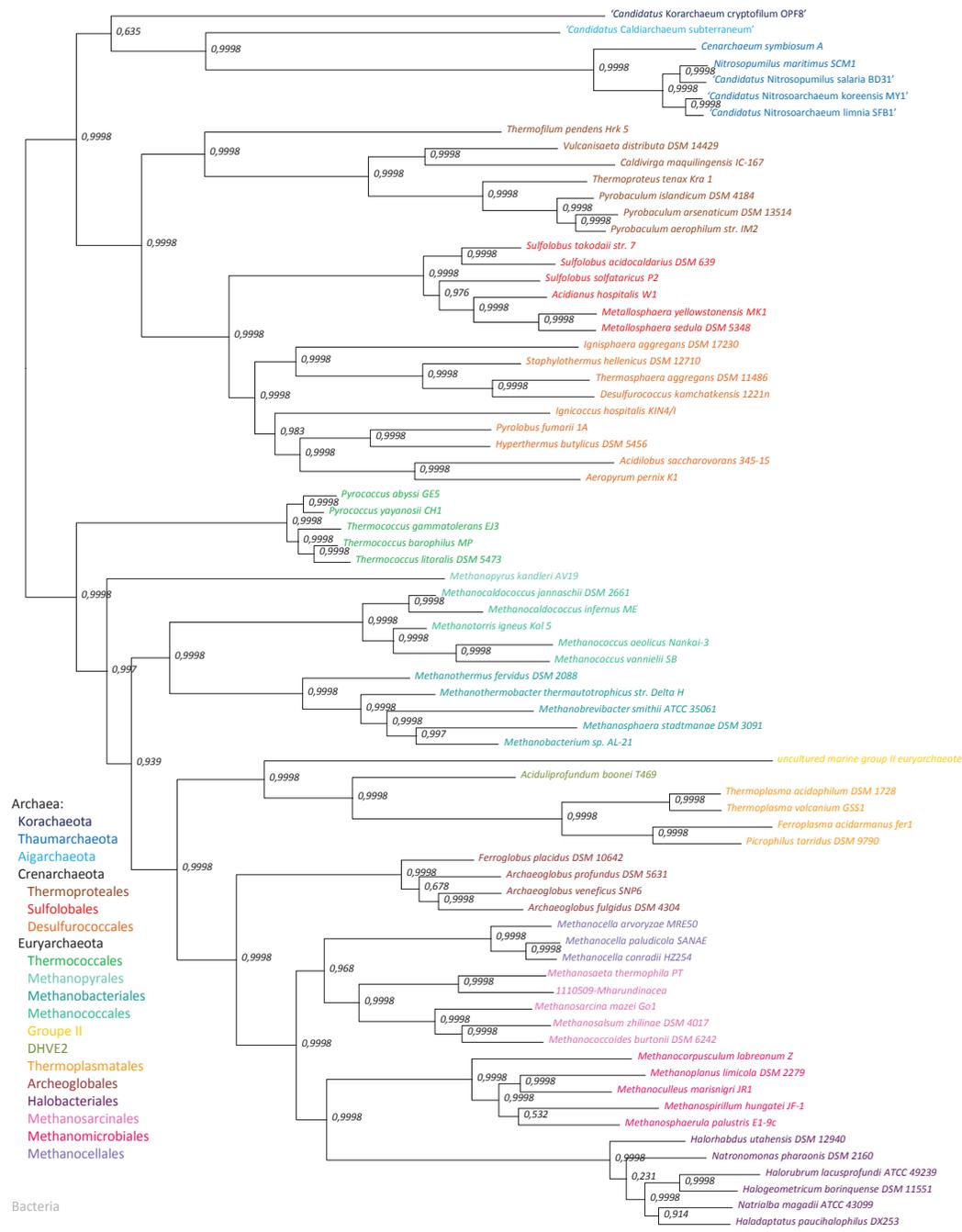
Petitjean et al.
 Supplementary Figure S5
 SF₁₉ PhyML-SH LG+G8



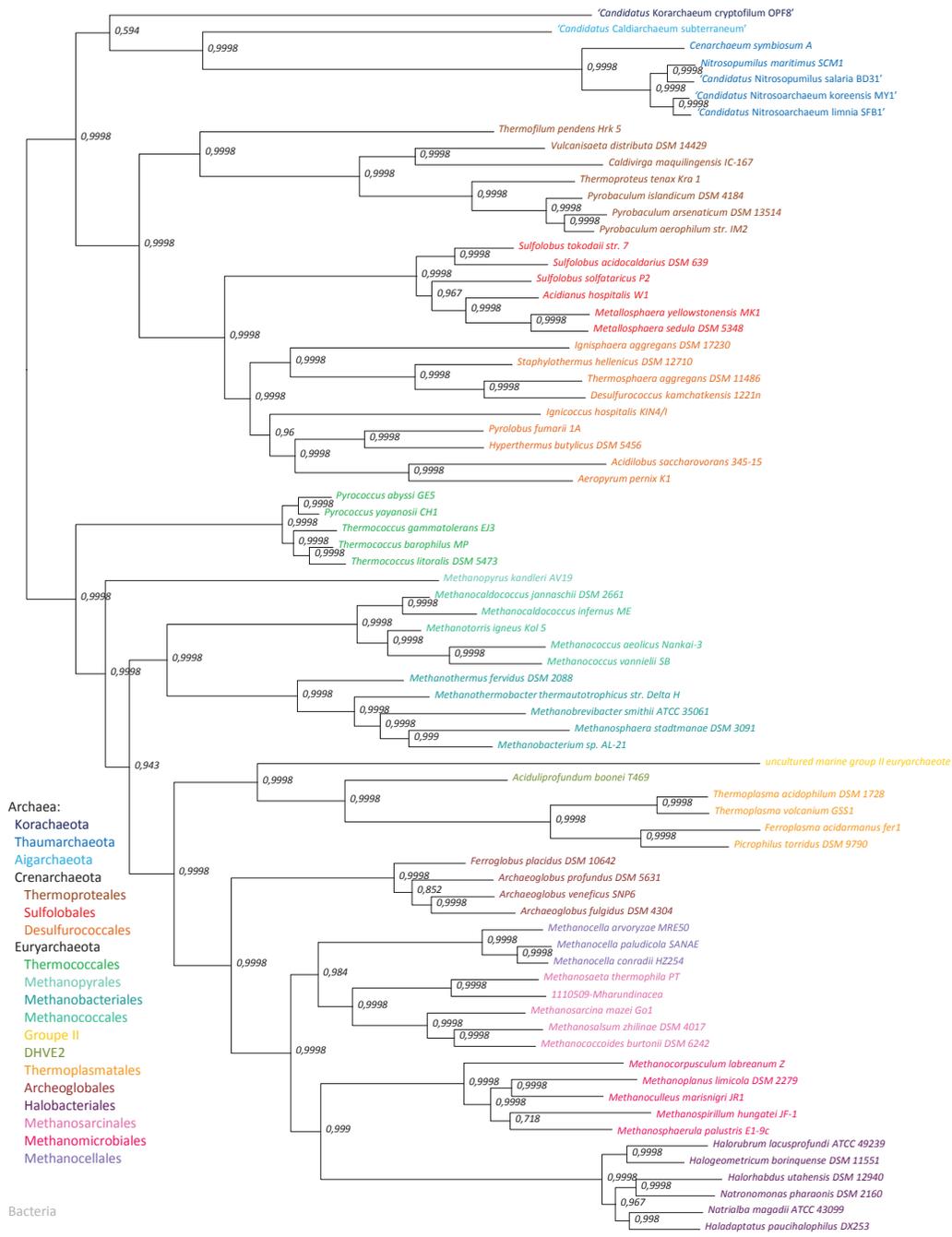
Petitjean et al.
 Supplementary Figure S5
 SF₂₀ PhyML-SH LG+G8



Petitjean et al.
Supplementary Figure S5
SF₂₂ PhyML-SH LG+G8



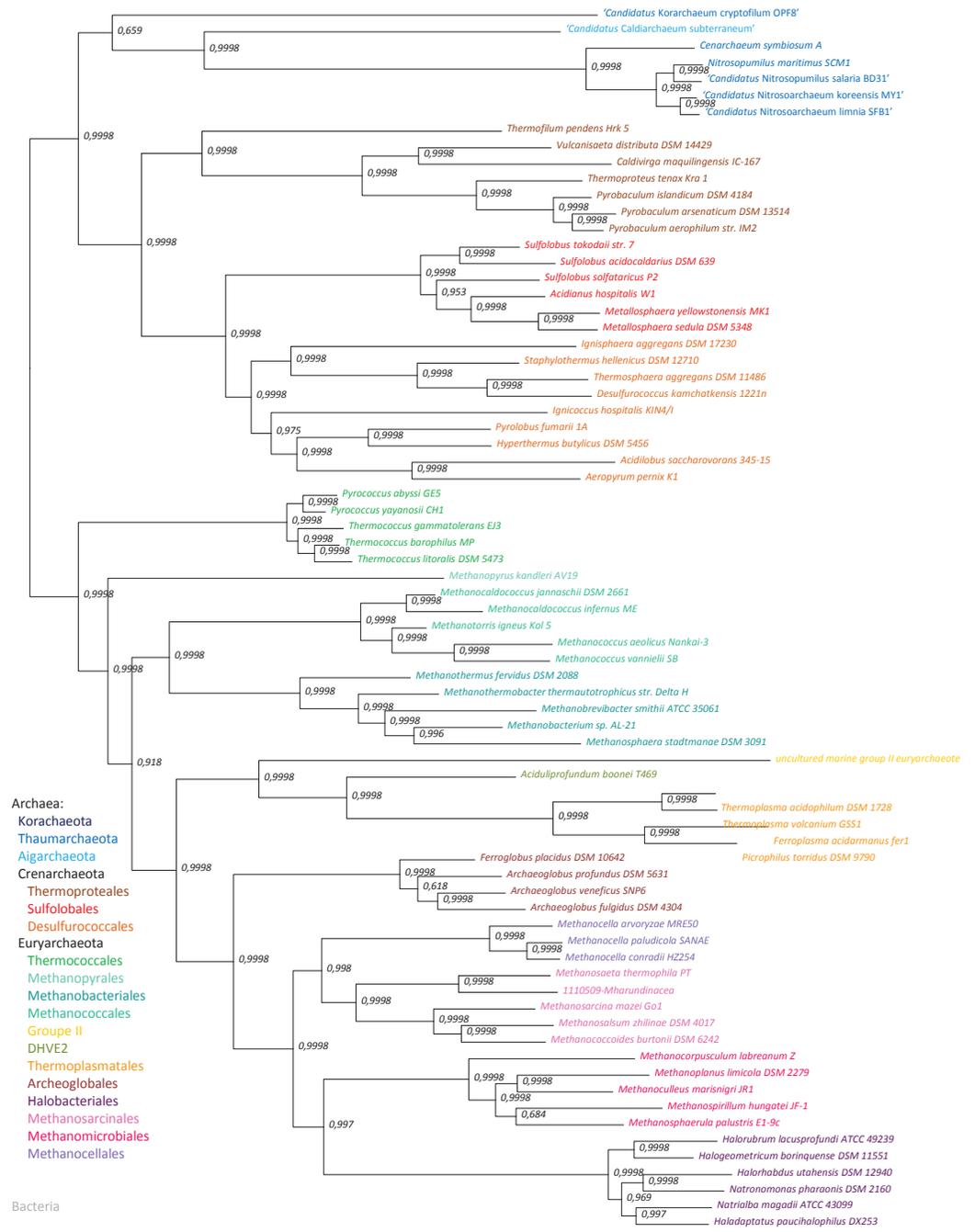
Petitjean et al.
Supplementary Figure S5
SF₂₃ PhyML-SH LG+G8



Petitjean et al.
Supplementary Figure S5
SF₂₄ PhyML-SH LG+G8

0.2

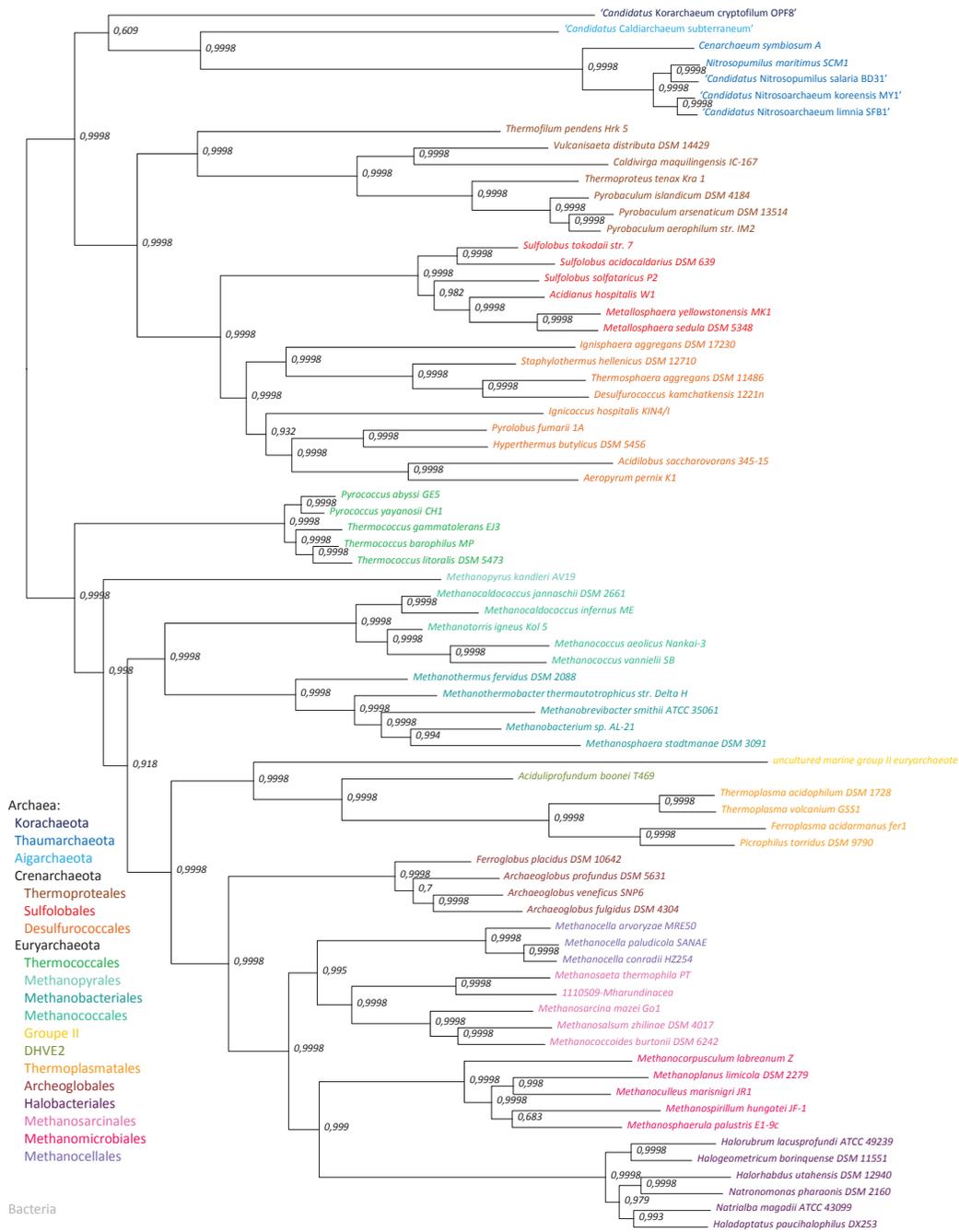
74 species – 31 225 positions



Petitjean et al.
Supplementary Figure S5
SF₂₆ PhyML-SH LG+G8

0.2

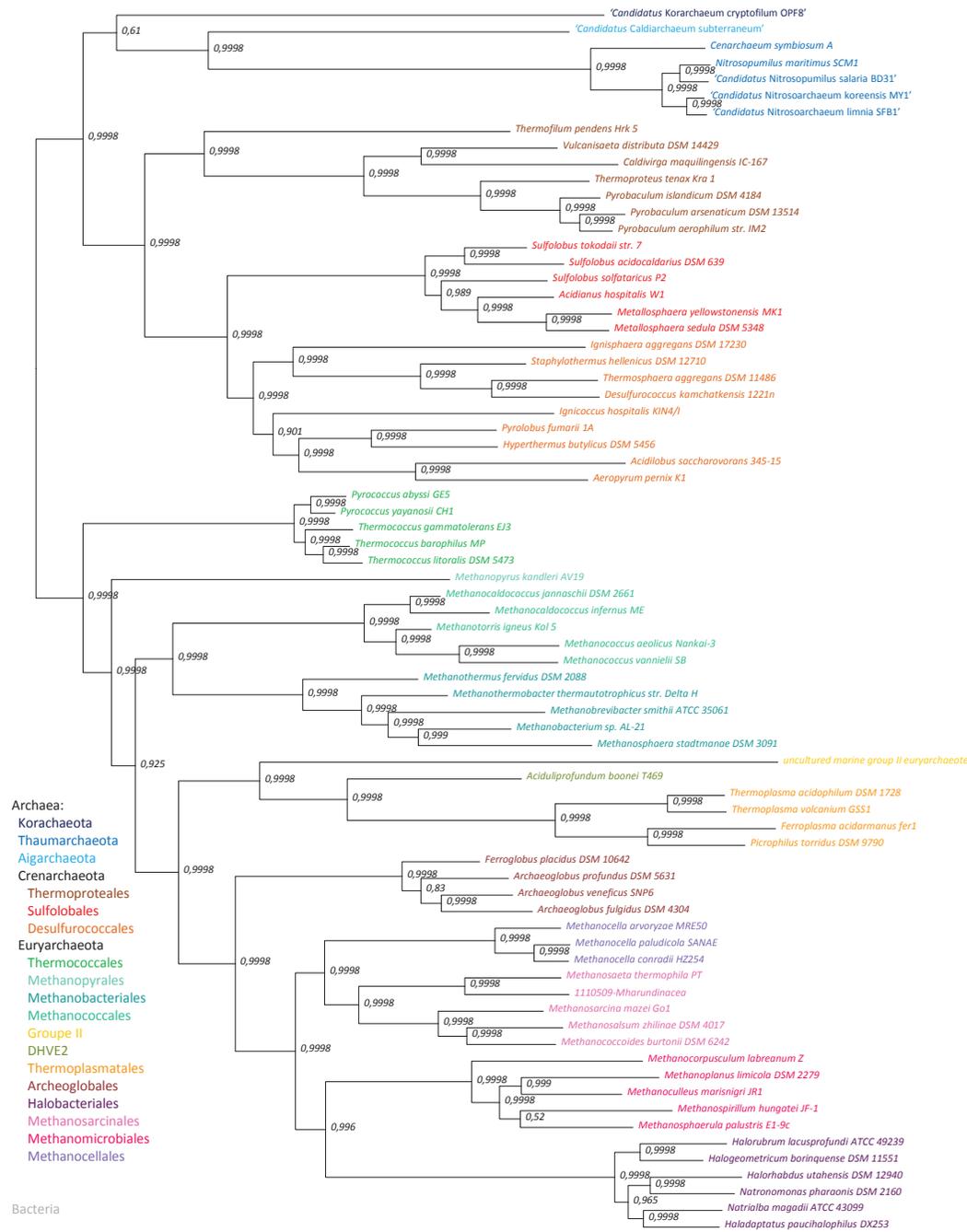
74 species – 22 472 positions



Petitjean et al.
 Supplementary Figure S5
 SF₂₇ PhyML-SH LG+G8

0.2

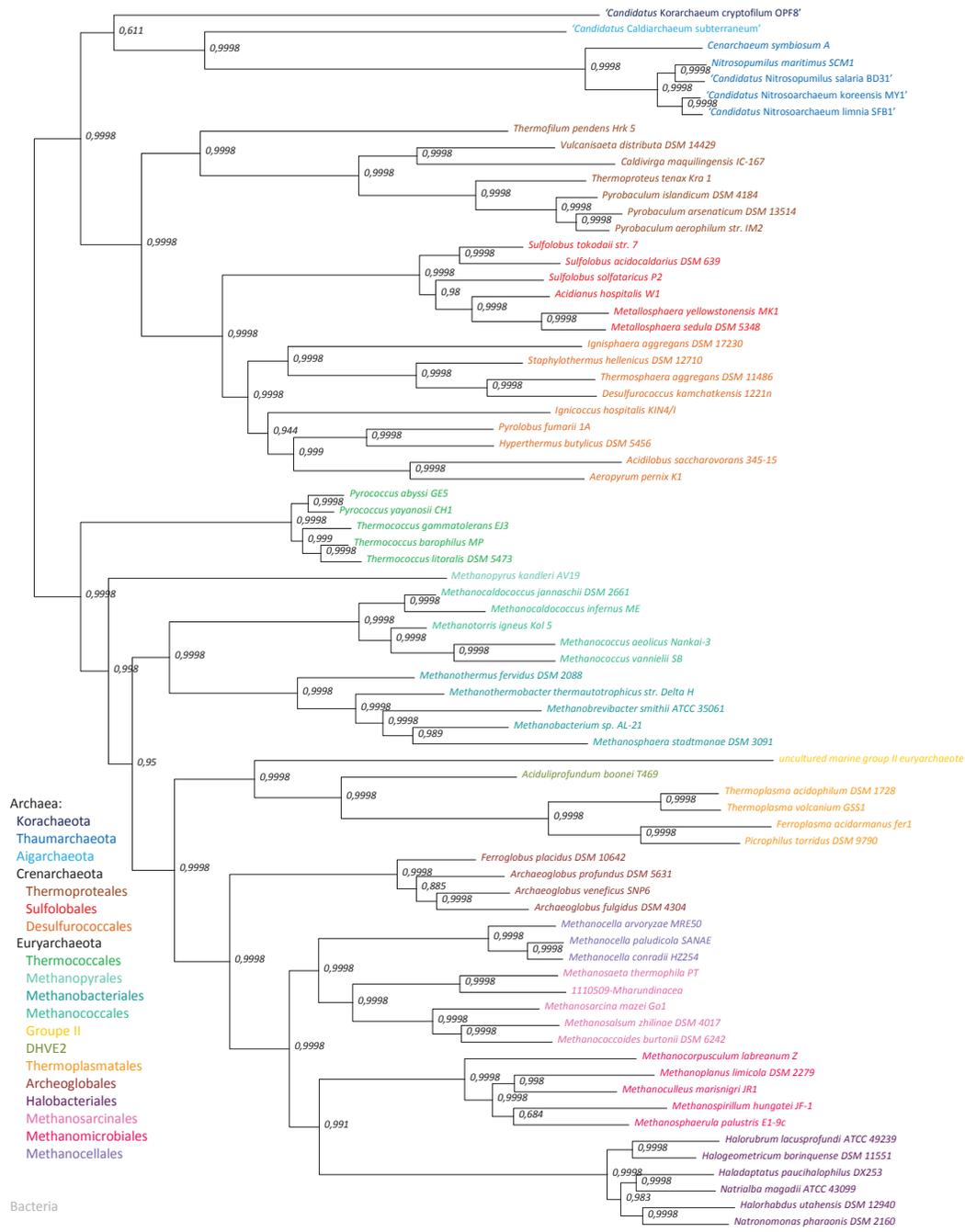
74 species – 32 999 positions



Petitjean et al.
 Supplementary Figure S5
 SF₂₉ PhyML-SH LG+G8

0.2

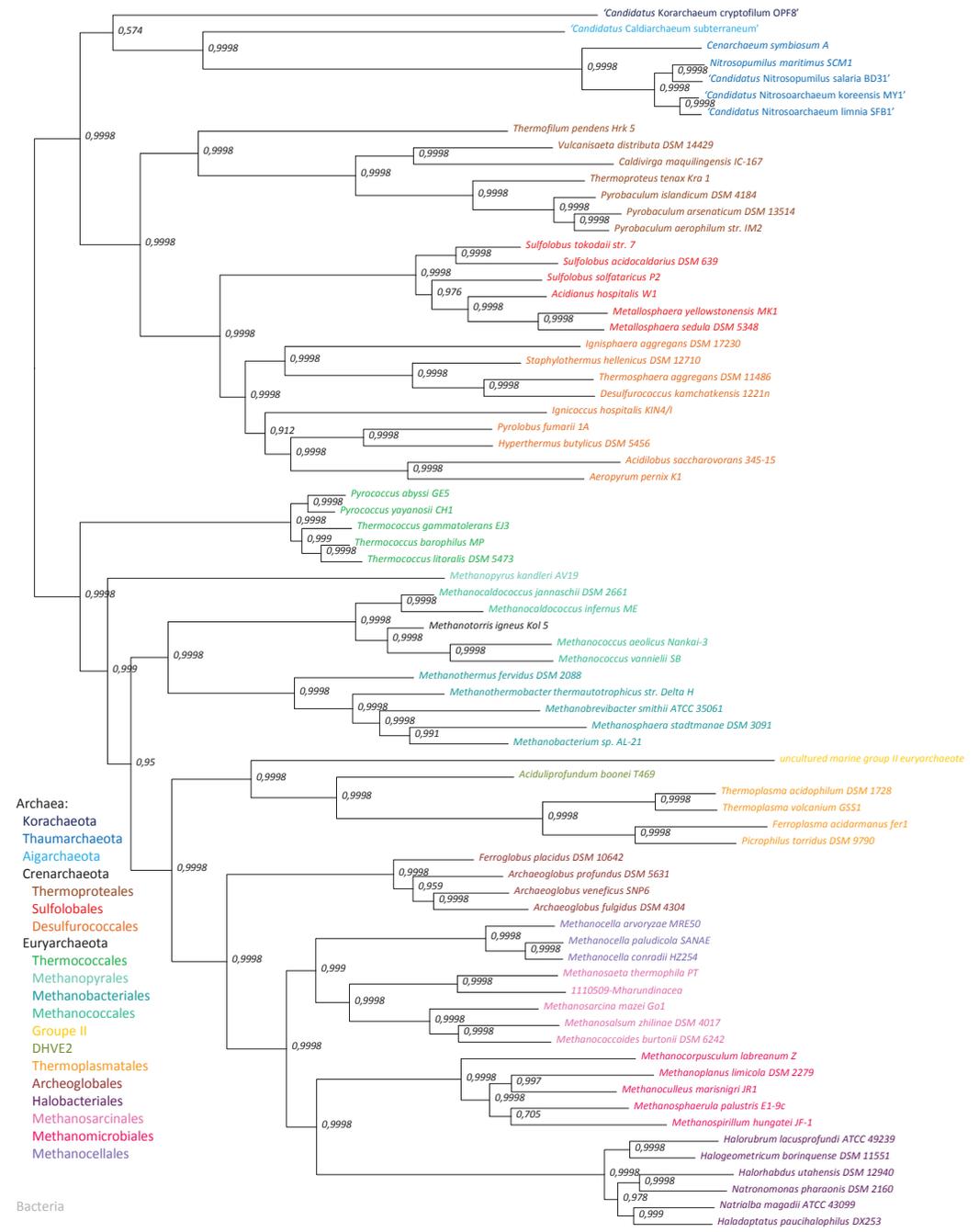
74 species – 33 799 positions



Petitjean et al.
Supplementary Figure S5
SF₃₁ PhyML-SH LG+G8

0.3

74 species – 34 449 positions



Petitjean et al.
Supplementary Figure S5
SF₃₄ PhyML-SH LG+G8

0.3

74 species – 35 589 positions

Annexe 2 : Listes et annotation des protéines des trois génomes étudiés.

- a. Liste des protéines de *Nitrosopumilus maritimus*
- b. Liste des protéines de *Cenarchaeum symbiosum*
- c. Liste des protéines de '*Candidatus Caldiarchaeum subterraneum*'

La colonne « Taille » correspond à la taille de la protéine en acides aminés.

Les colonnes « Alignement » correspondent aux nombre de positions et nombre de séquences des alignements générés automatiquement au départ de l'analyse.

Dans la colonne « Spéciale archées », une protéine annotée « 1 » correspond à l'absence d'homologue chez les bactéries et les eucaryotes, une protéine annotée « 0 » correspond à la présence d'au moins un homologue dans un de ces domaines.

Dans les autres colonnes « Archaea », « Eucarya » et « Bacteria », une protéine annotée « 1 » correspond à la présence d'au moins un homologue dans le groupe considéré, une protéine annotée « 0 » correspond à l'absence d'homologue dans ce groupe.

Dans la colonne « Marqueur », une protéine annotée « M », « M2 » ou « M3 » est considérée comme un marqueur potentiel de premier choix, deuxième choix ou troisième choix respectivement.

Dans la colonne « THG potentiel », sont notés les transferts horizontaux potentiels :

Les lettres correspondent à différents groupes taxonomiques :

A : Archées

B : Bactéries

E : Eucaryotes

T : Thaumarchées

EA : Euryarchées

CA : Crenarchées

K : Korarchées

AA : Aigarchées

Les flèches marquent le sens du transfert potentiel :

B=>A : Transfert potentiel de bactérie vers une ou plusieurs archées.

B<=>A : Transferts potentiel entre bactéries et archées, mais dont le sens reste à déterminer.

1349	161529271	YP_001583097.1	aminotransferase class V	381	244	200	0	1	1	1	1	1	1	1			B=>T
1350	161529272	YP_001583098.1	hypothetical protein Nmar_1764	198	192	5	1	0	0	0	0	0	0	0			
1351	161529273	YP_001583099.1	dFe-S ferredoxin iron-sulfur binding domain-containing protein	99	99	10	1	0	0	0	0	0	0	0			
1352	161529274	YP_001583100.1	adenylate/quarylate cyclase with integral membrane sensor	591	82	35	0	1	0	0	0	0	0	1			
1353	161529275	YP_001583101.1	heat shock protein DnaJ domain-containing protein	384	234	7	0	1	0	0	0	0	1	0			
1354	161529276	YP_001583102.1	UspA domain-containing protein	140	42	230	0	1	1	0	1	1	1	1			
1355	161529277	YP_001583103.1	amino acid permease-associated region	439	257	200	0	1	1	0	0	0	1	1			
1356	161529278	YP_001583104.1	hypothetical protein Nmar_1770	218	206	7	1	0	0	0	0	0	0	0			
1357	161529279	YP_001583105.1	glutamine synthetase type I	461	407	200	0	1	1	1	0	0	0	1			
1358	161529280	YP_001583106.1	lysine 2-scramonulase related protein	448	166	200	0	1	0	0	0	0	1	1			
1359	161529281	YP_001583107.1	MecS mechanosensitive ion channel	547	162	32	0	1	0	0	0	0	0	1			
1360	161529282	YP_001583108.1	ribose-phosphate pyrophosphokinase	292	217	200	0	1	1	1	0	1	1	1			M3
1361	161529283	YP_001583109.1	putative RNA methylase	315	99	196	0	1	1	1	1	1	1	0			
1362	161529284	YP_001583110.1	ribonuclease Z	299	137	200	0	1	1	1	1	1	1	1			M3
1363	161529285	YP_001583111.1	silent information regulator protein Sir2	242	101	200	0	1	1	0	1	1	1	1			
1364	161529287	YP_001583113.1	hypothetical protein Nmar_1779	134	88	43	1	1	1	1	0	0	0	0			M3
1365	161529288	YP_001583114.1	dephospho-CoA kinase GogE	197	101	105	0	1	1	0	0	1	0	0			
1366	161529289	YP_001583115.1	2p-5p RNA ligase	180	57	176	0	1	1	1	1	1	0	1			M3
1367	161529290	YP_001583116.1	HRNA cytidyltransferase	443	193	98	1	1	1	1	1	0	0	0			M
1368	161529291	YP_001583117.1	putative serine/threonine protein kinase	253	91	52	0	1	1	1	0	1	1	1			M3
1369	161529293	YP_001583119.1	FAD-dependent pyridine nucleotide-disulphide oxidoreductase	382	75	200	0	1	1	1	1	1	1	1			
1370	161529294	YP_001583120.1	short-chain dehydrogenase/reductase SDR	246	136	200	0	1	1	0	0	1	1	1			
1371	161529295	YP_001583121.1	ANI-type Zinc finger protein	73	73	10	1	1	0	0	0	0	0	0			
1372	161529296	YP_001583122.1	hypothetical protein Nmar_1788	94	89	4	1	0	0	0	0	0	0	0			
1373	161529297	YP_001583123.1	hypothetical protein Nmar_1789	185	126	5	1	0	0	1	0	0	0	0			
1374	161529298	YP_001583124.1	glutamine synthetase type I	490	409	200	0	1	1	1	0	0	0	1			
1375	161529300	YP_001583126.1	thymosin	570	461	200	1	1	1	1	1	0	0	0			M2
1376	161529301	YP_001583127.1	serine hydroxymethyltransferase	440	313	200	0	1	1	1	1	1	1	1			M2
1377	161529302	YP_001583128.1	DNA polymerase II large subunit	1125	900	75	1	1	0	1	1	1	0	0			M
1378	161529303	YP_001583129.1	geranylgeranylglyceryl phosphate synthase-like protein	250	174	75	0	1	1	1	1	0	1	1			M2
1379	161529304	YP_001583130.1	hypothetical protein Nmar_1796	120	119	4	1	0	0	0	0	0	0	0			
1380	161529305	YP_001583131.1	hypothetical protein Nmar_1797	200	123	6	1	0	0	1	0	0	0	0			
1381	161529307	YP_001583133.1	hypothetical protein Nmar_1799	128	125	5	1	0	0	1	0	0	0	0			

578	118577124	YP_876867.1	ABC-type multidrug transport system ATPase component	313	168	200	0	1	1	1	1	0	0	0	
579	118577126	YP_876869.1	hypothetical protein CENSyA_1958	347	187	49	0	1	1	1	1	0	1	M2	A=>B
580	118577128	YP_876871.1	molecular chaperone GrpE	187	85	114	0	1	0	0	0	1	1		
581	118577130	YP_876873.1	DnaJ-class molecular chaperone	351	278	200	0	1	0	0	0	1	1		
582	118577134	YP_876877.1	archaeal Glu-IRNAGln amidotransferase subunit E	628	450	106	0	1	1	1	1	0	1	M	
583	118577135	YP_876878.1	archaeal Glu-IRNAGln amidotransferase subunit D/asparaginase	392	170	200	0	1	1	1	1	1	1	M	
584	118577136	YP_876879.1	AAA ATPase	728	568	200	0	1	1	1	1	1	1	M	
585	118577141	YP_876884.1	hypothetical protein CENSyA_1975	107	128	6	1	0	0	0	0	0	0		
586	118577143	YP_876886.1	dimucellin-binding enzyme	206	158	61	0	1	0	1	0	0	1	M2	E=>B
587	118577144	YP_876887.1	fructose-2,6-bisphosphatase	228	42	200	0	1	1	0	0	1	1		
588	118577152	YP_876895.1	thiol-disulfide isomerase	135	56	200	0	1	1	1	1	0	1		
589	118577153	YP_876896.1	hypothetical protein CENSyA_1988	139	77	82	0	1	1	1	1	1	1	M2	E=>B
590	118577156	YP_876899.1	polyphenyltransferase	310	203	200	0	1	1	1	1	1	0	M3	
591	118577157	YP_876900.1	transcriptional regulator	152	53	200	0	1	1	1	1	0	0		
592	118577159	YP_876902.1	threonine-phosphate decarboxylase	395	113	200	0	1	1	1	0	1	1		
593	118577160	YP_876903.1	cobyrinic acid synthase	276	287	200	0	1	0	0	0	0	0	1	
594	118577161	YP_876904.1	cobalamin biosynthesis protein	320	122	200	0	1	1	0	0	1	1	M3	A=>B/E
595	118577162	YP_876905.1	cobalamin-5-phosphate synthase	219	175	35	0	1	0	0	0	0	1	M3	
596	118577163	YP_876906.1	4-diphosphocytidyl-2C-methyl-D-erythritol synthase	198	99	54	0	1	1	0	0	0	0	M2	
597	118577165	YP_876908.1	adenine deaminase	609	290	133	0	1	1	1	1	1	1		
598	118577167	YP_876910.1	ROK-family glucokinase	282	122	200	0	1	1	0	0	0	1		B=>T
599	118577169	YP_876912.1	hypothetical protein CENSyA_2004	207	141	12	1	0	1	1	0	0	0	M3	
600	118577171	YP_876914.1	uncharacterized Rossmann fold enzyme	239	104	89	1	1	1	1	1	0	0	M2	
601	118577172	YP_876915.1	dihydropterate synthase	275	136	200	0	1	1	0	1	1	1		
602	118577179	YP_876922.4	hypothetical protein CENSyA_2014	362	162	9	1	1	1	0	0	0	0		
603	118577181	YP_876924.1	Zn-dependent oxidoreductase	442	151	200	0	1	1	1	1	1	1		T<=>B
604	118577182	YP_876925.1	TPR repeat protein	206	45	178	0	1	0	0	0	1	1		
605	118577188	YP_876931.1	DNA-binding protein containing a Zn-ribbon domain	508	202	98	0	1	1	1	0	0	0		M
606	118577189	YP_876932.1	glucose/sorbitose dehydrogenase	390	114	200	0	1	1	1	0	1	0		
607	118577190	YP_876933.1	hydroxymethylglutaryl-CoA reductase	415	288	69	0	1	1	1	1	1	1		
608	118577191	YP_876934.1	hypothetical protein CENSyA_2027	191	129	15	0	1	0	0	0	0	1		T<=>EA/B
609	118577192	YP_876935.1	hypothetical protein CENSyA_2028	193	142	15	1	1	0	0	0	0	0		T<=>EA
610	118577193	YP_876936.1	hypothetical protein CENSyA_2029	208	142	15	1	1	0	0	0	0	0		T<=>EA
611	118577194	YP_876937.1	hypothetical protein CENSyA_2030	218	145	11	1	0	0	0	0	0	0		
612	118577195	YP_876938.1	hypothetical protein CENSyA_2031	183	122	16	0	1	0	0	0	0	1		T<=>EA/B
613	118577196	YP_876939.1	hypothetical protein CENSyA_2032	207	142	15	1	1	0	0	0	0	0		T<=>EA
614	118577197	YP_876940.1	hypothetical protein CENSyA_2033	189	142	15	1	1	0	0	0	0	0		T<=>EA
615	118577200	YP_876943.1	histidinol dehydrogenase	422	265	200	0	1	1	1	1	0	1		A=>B
616	118577201	YP_876944.1	histidinol-phosphate aminotransferase	354	120	200	0	1	1	1	0	1	1		
617	118577202	YP_876945.1	phosphatase	321	258	17	1	0	1	0	0	0	0		
618	118577204	YP_876947.1	glutamine amidotransferase	201	116	200	0	1	1	1	0	1	1		
619	118577205	YP_876948.1	phosphoribosyl formimino-5-aminoimidazole isomerase	234	149	200	0	1	0	1	0	0	1		B=>AA+T
620	118577210	YP_876953.1	4-methyl-5-(beta-hydroxyethyl)thiazole monophosphate synthesis protein	458	169	200	0	1	0	0	1	1	1		
621	118577215	YP_876958.1	hypothetical protein CENSyA_2051	204	119	16	0	1	0	0	0	0	1		T<=>EA/B
622	118577223	YP_876966.1	hypothetical protein CENSyA_2059	1846	93	24	1	0	1	0	0	0	0		
623	118577224	YP_876967.1	hypothetical protein CENSyA_2060	1232	167	22	1	0	1	0	0	0	0		
624	118577225	YP_876968.1	hypothetical protein CENSyA_2061	473	911	10	1	0	0	0	0	0	0		
625	118577228	YP_876971.1	hypothetical protein CENSyA_2064	523	84	9	1	0	0	0	0	0	0		
626	118577230	YP_876973.1	surface layer-associated STABLE protease	1047	122	110	0	1	1	0	0	1	1		

1061	315427598	BAJ49197.1	UDP-glucose 4-epimerase	315	173	100	0	1	1	0	1	0	1
1062	315427603	BAJ49202.1	NAD-dependent epimerasedehydratase	347	150	100	0	1	1	1	1	0	1
1063	315427739	BAJ49334.1	conserved hypothetical protein	87	194	12	1	0	1	0	1	0	0
1064	315427751	BAJ49346.1	hypothetical protein HGMM_F01D06C13	695	557	6	0	1	0	0	0	0	1
1065	315427767	BAJ49362.1	arsenic transporting ATPase	327	131	100	0	1	1	0	0	1	1
1066	315427770	BAJ49365.1	carbon starvation protein CstA	596	376	100	0	1	1	0	0	0	1
1067	315427781	BAJ49375.1	enoyl-CoA hydratase	254	192	100	0	1	1	0	1	1	1
1068	315427782	BAJ49376.1	alcohol dehydrogenase	339	178	100	0	1	1	0	0	1	1
1069	315427783	BAJ49377.1	naphthoate synthase	312	220	100	0	1	1	0	0	1	1
1070	315427784	BAJ49378.1	aldehydeferredoxin oxidoreductase	627	373	100	0	1	1	0	1	0	1
1071	315427786	BAJ49380.1	phenylacetate-CoA ligase	458	375	100	0	1	1	0	0	0	1
1072	315427790	BAJ49384.1	branched-chain amino acid ABC transporter substrate-binding protein	494	222	22	0	1	0	0	0	0	1
1073	315427795	BAJ49389.1	phenylacetyl-CoA acceptor oxidoreductase subunit	211	132	100	0	1	1	0	0	1	1
1074	315427796	BAJ49390.1	phenylacetyl-CoA acceptor oxidoreductase subunit	904	199	100	0	1	1	0	0	0	1
1075	315427797	BAJ49391.1	benzoyl-CoA reductase subunit C	394	205	35	0	1	0	0	0	0	1
1076	315427798	BAJ49392.1	CoA-substrate-specific enzyme activase	249	137	100	0	1	0	0	0	0	1
1077	315427799	BAJ49393.1	CoA-substrate-specific enzyme activase	275	138	100	0	1	0	0	0	0	1
1078	315427800	BAJ49394.1	benzoyl-CoA reductase subunit B	420	171	29	0	1	0	0	0	0	1
1079	315427969	BAJ49559.1	dihydroorotate oxidase	332	207	100	0	1	0	0	0	1	1
1080	315427983	BAJ49572.1	5-oxoprolinase ATP-hydrolysing	498	399	100	0	1	1	0	0	0	1
1081	315427984	BAJ49573.1	N-methylhydantoinase B	566	304	100	0	1	1	0	0	0	1
1082	315427986	BAJ49575.1	N-methylhydantoinase A	692	398	100	0	1	1	0	0	0	1
1083	315427989	BAJ49578.1	cobaltnickel ABC transporter ATP-binding protein	586	236	100	0	1	1	0	1	0	1

Annexe 3 : Répartition taxonomique des 200 nouveaux marqueurs chez les archées.

Chaque ligne correspond à une espèce, chaque colonne à un jeu de données, le croisement entre les deux correspond à la présence (« 1 », case verte), ou à l'absence (« 0 », case rouge) d'homologue de l'espèce dans le jeu de données. La dernière colonne donne le nombre de jeux de données dans lesquels l'espèce est présente.

La liste des archées correspond uniquement aux 129 génomes utilisés dans ce travail.

Annexe 4 : Liste des protéines informationnelles mises à jour pour le Chapitre 1.

Liste des 57 protéines ribosomiques, 14 sous-unités de l'ARN polymérase et 2 facteurs de transcription mis à jours pour l'étude exposée dans le Chapitre 1.

Annotation	Référence gi	Référence accession	séquences
50S ribosomal protein L1	161527890	YP_001581716	127
50S ribosomal protein L2	161527614	YP_001581440	128
50S ribosomal protein L3P	161528317	YP_001582143	127
50S ribosomal protein L4P	161528316	YP_001582142	127
50S ribosomal protein L5	161528305	YP_001582131	128
50S ribosomal protein L6	161528302	YP_001582128	128
50S ribosomal protein L7Ae	161527733	YP_001581559	129
acidic ribosomal protein P0 - L10	161527889	YP_001581715	128
50S ribosomal protein L10e	161527957	YP_001581783	129
50S ribosomal protein L11	161527893	YP_001581719	127
50S ribosomal protein L12 - L12e	161527882	YP_001581708	128
50S ribosomal protein L13	161527934	YP_001581760	129
50S ribosomal protein L14	161528308	YP_001582134	127
50S ribosomal protein L15	161527909	YP_001581735	129
50S ribosomal protein L15e	161528762	YP_001582588	129
50S ribosomal protein L18	161527906	YP_001581732	128
50S ribosomal protein L15 – L18e	161527935	YP_001581761	127
50S ribosomal protein L19e	161527898	YP_001581724	128
50S ribosomal protein L21e	161528892	YP_001582718	128
50S ribosomal protein L22	161528313	YP_001582139	128
50S ribosomal protein L25 -L23	161528315	YP_001582141	127
KOW domain-containing protein L24	161528307	YP_001582133	128
50S ribosomal protein L24 - L24e	161527735	YP_001581561	127
50S ribosomal protein L29	161528311	YP_001582137	125
50S ribosomal protein L30	161527908	YP_001581734	129
50S ribosomal protein L31e	161528771	YP_001582597	126
50S ribosomal protein L32e	161527899	YP_001581725	128
hypothetical protein Nmar_0434 -L37ae	161527942	YP_001581768	129
50S ribosomal protein L37e	161527728	YP_001581554	122
50S ribosomal protein L39e	161528772	YP_001582598	121
50S ribosomal protein L40e	161528633	YP_001582459	126
50S ribosomal protein L44e	161529087	YP_001582913	127
30S ribosomal protein S2	161527824	YP_001581650	129
30S ribosomal protein S3	161528312	YP_001582138	128
30S ribosomal protein S3Ae	161529020	YP_001582846	128
30S ribosomal protein S4	161527832	YP_001581658	129
30S ribosomal protein S4e	161528306	YP_001582132	128
30S ribosomal protein S5	161527907	YP_001581733	129
30S ribosomal protein S6e	161527583	YP_001581409	129
30S ribosomal protein S7	161527863	YP_001581689	129
30S ribosomal protein S8	161528303	YP_001582129	128
30S ribosomal protein S8e	161528028	YP_001581854	129
30S ribosomal protein S9	161527933	YP_001581759	129
30S ribosomal protein S10	161528541	YP_001582367	129
30S ribosomal protein S11	161528958	YP_001582784	129
30S ribosomal protein S12	161527862	YP_001581688	129
30S ribosomal protein S13	161527833	YP_001581659	129
30S ribosomal protein S14	161528304	YP_001582130	125
30S ribosomal protein S15	161529016	YP_001582842	128
30S ribosomal protein S17	161528309	YP_001582135	127
hypothetical protein Nmar_1430 -S17e	161528938	YP_001582764	128
30S ribosomal protein S19	161528314	YP_001582140	128
30S ribosomal protein S19e	161529132	YP_001582958	129
hypothetical protein Nmar_0535 – S24e	161528043	YP_001581869	127
30S ribosomal protein S27ae	161528042	YP_001581868	123
30S ribosomal protein S27E	161529088	YP_001582914	127
30S ribosomal protein S28e	161527734	YP_001581560	129
RNA polymerase Rbp10 – RpoP	161528540	YP_001582366	126
DNA-directed RNA polymerase subunit N – RpoN	161527826	YP_001581652	129
transcription termination factor Tfs – RpoM	161529264	YP_001583090	125
hypothetical protein Nmar_1757 – RpoL	161529265	YP_001583091	127
RNA polymerase Rpb6 – RpoK	161528440	YP_001582266	127
RNA polymerase Rpb5 – RpoH	161527854	YP_001581680	128
RNA polymerase Rpb4 – RpoF	161528891	YP_001582717	127
DNA-directed RNA polymerase subunit E RpoE2	161528744	YP_001582570	127
DNA-directed RNA polymerase subunit E' – RpoE1	161528743	YP_001582569	129
RNA polymerase insert RpoD	161527936	YP_001581762	129
DNA-directed RNA polymerase subunit B – RpoB2	161527855	YP_001581681	128
DNA-directed RNA polymerase subunit B – RpoB1	161527856	YP_001581681	128
bifunctional DNA-directed RNA polymerase subunit A'/A" – RpoA2	161527856	YP_001581682	128
bifunctional DNA-directed RNA polymerase subunit A'/A" – RpoA1	161527857	YP_001581683	129
NusG antitermination factor	161527894	YP_001581720	127
NusA family protein	161527861	YP_001581687	127

Annexe 5 : Liste des matrices utilisées pour l'analyse en désaturation par sélection de sites pour l'étude exposée dans le Chapitre 1.

Matrice	Nombre de positions
S ₀	1958
S ₁	3884
S ₂	5439
S ₃	6929
S ₄	8373
S ₅	9742
S ₆	11089
S ₇	12465
S ₈	13806
S ₉	15176
S ₁₀	16508
S ₁₁	17845
S ₁₂	19167
S ₁₃	20458
S ₁₄	21746
S ₁₅	23000
S ₁₆	24143
S ₁₇	25237
S ₁₉	27217
S ₂₀	28146
S ₂₂	29825
S ₂₃	30561
S ₂₄	31225
S ₂₆	32472
S ₂₇	32999
S ₂₉	33799
S ₃₁	34449
S ₃₄	35102
Totale	35589

Annexe 6 : Liste des matrices utilisées pour l'analyse en désaturation par sélection de gènes pour l'étude exposée dans le Chapitre 1.

Pour chaque marqueur, la matrice indiquée est la plus grande dans laquelle ce marqueur a été inclus.

Par exemple, les marqueurs de la matrice DT30 sont aussi dans les matrices DT20 et DT10.

Le nombre de distance calculées correspond au nombre de couples d'espèces pour lesquels une distance a été mesurée, et qui sert de diviseur pour le calcul de la distance moyenne (maximum = 10).

Marqueur	Distance moyenne	Nombre de distances calculées	Matrice	Nombre de positions
s28e	0,81314591	10	DT10	1242
s19	0,87391316	10		
l12e	0,93431001	10		
rpoM	0,97764542	10		
l37e	0,99762814	10		
M036	1,05088415	10		
M046	1,05354258	10		
l7ae	1,0705633	10		
M066	1,07546077	6		
M197	1,12731502	6		
MA09	1,14941092	6	DT20	2672
rpoK	1,15097537	10		
s12	1,15257442	10		
s2	1,16728252	10		
s7	1,18210045	10		
s11	1,18845449	10		
l39e	1,19038229	10		
l5	1,2145064	10		
M128	1,24681426	10		
M039	1,25275645	10		
MA23	1,27057753	10	DT30	4701
M086	1,27195784	10		
M201	1,29563716	10		
M207	1,29925848	10		
s19e	1,32980492	10		
M171	1,33274006	10		
M136	1,33417573	10		
MA55	1,34222622	10		
l3	1,34508504	10		
rpoN	1,35059768	10		
M070	1,38272673	10	DT40	6330
MA42	1,39090948	6		
M147	1,39397627	10		
M194	1,41037915	10		
MA44	1,41982895	10		
s5	1,42965314	10		
M159	1,43083271	10		
nusA	1,43475249	10		
M122	1,44850164	10		

M053	1,45048367	10	DT50	8733
MA22	1,45300464	10		
M187	1,45823695	10		
s15	1,45895579	10		
MA52	1,45955437	6		
MA56	1,46319707	10		
M138	1,46852341	10		
M219	1,47410339	10		
rpoA1	1,47653898	10		
l15e	1,4770368	10		
M173	1,48723625	10		
M213	1,48769448	6		
MA54	1,48959821	10		
rpoB2	1,49780334	10		
s14	1,50123284	10		
s6e	1,50692094	10		
M184	1,5087578	10		
MA02	1,52833851	10		
l6	1,5326468	10		
M148	1,53295574	10		
M215	1,53507611	10		
l18	1,53508518	10		
M210	1,54374317	6		
l21e	1,54800109	10		
M126	1,55528883	10		
M037	1,55583385	10		
MA33	1,55617389	10		
s9	1,56276877	10		
l40e	1,57209471	10		
s17e	1,57961288	10		
s27ae	1,58240479	10		
M212	1,58492593	10		
M182	1,58702595	6		
M164	1,58786794	10		
M178	1,59166906	10		
MA20	1,59818404	10		
rpoB1	1,6056481	10		
M134	1,60834848	10		
M028	1,60895551	10		
M125	1,60974439	10		
s8e	1,61139908	10		
MA37	1,61157562	10		
MA08	1,61504253	10		
s10	1,62061821	10		
l11	1,62389996	10		
s27e	1,62453301	10		
M190	1,6253055	10		
s8	1,62783873	10		
M032	1,63652623	3		
M040	1,6390425	10		
			DT60	10963
			DT70	12840
			DT80	16340
			DT90	18552

M175	1,64378792	6	DT100	20629
M048	1,64612471	10		
l19e	1,64926715	10		
MA48	1,6518101	10		
s17	1,66160866	10		
M112	1,66920497	10		
M014	1,67635768	10		
MA53	1,67638445	10		
M145	1,67818156	10		
M115	1,68043672	10		
M193	1,68251642	10		
rpoE2	1,6835943	10	DT110	22987
M199	1,68750593	3		
rpoH	1,69432393	10		
M150	1,69511077	10		
M067	1,69860747	6		
M054	1,70618695	10		
l2	1,70618973	10		
M047	1,70769043	10		
M183	1,71158673	10		
M179	1,71368038	6		
M132	1,7166913	10		
M038	1,72678923	6		
M055	1,72787748	10		
M149	1,72921748	10		
M203	1,73160633	3		
M137	1,73293096	10		
s3ae	1,73610174	10		
MA05	1,73957769	10		
l22	1,74150042	10		
M043	1,74290963	10		
M205	1,75559127	10	DT130	27863
M002	1,75609693	10		
MA21	1,7635065	10		
MA34	1,76999781	10		
M166	1,77032753	10		
M069	1,77593921	10		
M160	1,77648761	10		
M198	1,77753262	10		
MA47	1,78662984	10		
M154	1,78957167	10		
l24e	1,79649302	10		
l24	1,79701265	10		
MA32	1,8027434	10		
l32e	1,82365418	10		
M218	1,82441005	10		
M140	1,82609143	10		
M121	1,83672235	10		
rpoA2	1,83942755	10		
M170	1,83986158	10		

M010	1,84349775	10	DT150	32350
M144	1,84549012	10		
M119	1,84695156	10		
M033	1,84957126	10		
M050	1,85102774	10		
MA06	1,85697567	3		
M189	1,8603604	6		
M041	1,86730981	10		
rpoD	1,87969418	10		
MA27	1,88091723	6	DT160	34383
M180	1,88631171	10		
I14	1,88733091	10		
M042	1,89380325	10		
I4	1,89637119	10		
M176	1,89869073	6		
MA17	1,90345972	6		
I1	1,915029	10		
M064	1,9263975	3		
MA10	1,9276212	1		
I37ae	1,92827006	10		
M045	1,94010169	10		
M151	1,94109423	10		
I10e	1,94242958	10		
M031	1,94342384	10		
M158	1,94502309	10		
M018	1,95293228	10		
M118	1,95612123	3	DT180	40222
M056	1,97297299	10		
M044	1,97307409	10		
MA25	1,97636206	10		
M029	1,9846848	10		
MA11	1,9968463	3		
M163	2,00887608	10		
M188	2,01565706	10		
I23	2,0258275	10		
M208	2,0295969	1	DT190	41985
M085	2,02974738	10		
s3	2,03084497	10		
M011	2,03361637	10		
MA01	2,03673252	10		
M186	2,03939054	10		
M021	2,05283813	10		
M135	2,05393557	10		
M153	2,06065383	10		
M143	2,06608771	10		
M206	2,06818987	10		
rpoE1	2,06850103	10		
M076	2,06999741	10		
M004	2,0786757	1		
MA03	2,07922682	10		

M027	2,08069095	10	DT200	44233
M057	2,09163333	10		
M181	2,0924164	10		
M078	2,09626543	10		
M129	2,0979118	3		
MA41	2,09916525	10		
M080	2,10355579	10		
M165	2,10695921	10	DT210	47068
M058	2,10953107	10		
M093	2,11022341	10		
M195	2,11249604	10		
M077	2,11973943	6		
M071	2,12595264	10		
MA30	2,13438893	10		
M035	2,13569968	10		
M025	2,13917813	10		
M030	2,14694397	10		
s4e	2,15617214	10	DT220	48803
M061	2,16334797	10		
MA36	2,16627175	10		
l30	2,17080781	10		
M052	2,1814731	10		
M142	2,18244907	6		
M204	2,18913913	6		
M117	2,19676109	10		
M162	2,22253408	10		
M156	2,22317808	10		
l31e	2,23807342	10	DT230	50694
MA45	2,24108395	10		
M123	2,24461605	10		
l18e	2,25657231	10		
M168	2,26454572	10		
M177	2,26865965	10		
l13	2,30213907	10		
M185	2,30378868	10		
M026	2,30489167	10		
rpoF	2,30735296	10		
M191	2,31472467	10	DT240	52179
nusG	2,31501369	10		
MA35	2,330827	10		
M192	2,33192669	10		
MA31	2,33737627	3		
l10	2,33958011	10		
M012	2,34153314	10		
l29	2,34431401	10		
M049	2,35663201	10		
MA50	2,36046493	10		
M209	2,38441758	6		
MA46	2,38842761	10		
M005	2,403694	1		

M051	2,42672475	10		
M124	2,45930511	10		
M161	2,46849158	10		
M157	2,47051129	10		
l15	2,47296853	10		
MA19	2,50098472	10	DT260	55565
MA18	2,5013873	3		
M127	2,5057969	3		
M060	2,51569924	10		
M202	2,51595733	10		
l44e	2,51782973	10		
M075	2,55022911	10		
MA40	2,58724886	10		
MA28	2,58987802	6		
M022	2,5985988	10		
M116	2,60849567	10		
M196	2,67002588	10		
M003	2,67553257	10		
M059	2,68817926	10		
rpoP	2,7186661	10		
MA39	2,7404339	10		
M146	2,75250332	10		
M130	2,79218596	10		
M092	2,82906282	10		
M174	2,83338458	10		
M034	2,84542223	10		
M217	2,87419791	10		
M200	2,93473747	3		
MA16	2,9804462	10		
s24e	3,03426218	10		
M068	3,08430638	10		
rpoL	3,55880201	10		
			DT270	57148
			DT280	58433

Annexe 7 : Matériels supplémentaires de l'Article 2 (Chapitre 2).

Supplementary material

Supplementary table S1. List of archaeal and bacterial complete genome sequences used in this work.

Supplementary table S2. List of all protein markers used in this work.

Supplementary table S3. List of species-pairs used for evolutionary estimations for gene-by-gene desaturation analysis.

Supplementary fig. S1. Maximum likelihood tree based on the complete set of markers (32 ribosomal proteins and 38 diverse conserved proteins, 9450 sites) and rooted using bacterial sequences as outgroup. Numbers at branches are bootstrap proportions. The scale bar indicates the number of substitutions per position.

Supplementary fig. S2. Unrooted maximum likelihood tree based on the complete set of markers (32 ribosomal proteins and 38 diverse conserved proteins, 9450 sites). Numbers at branches are bootstrap proportions. The scale bar indicates the number of substitutions per position.

Supplementary figs. S3-S12. Rooted maximum likelihood trees corresponding to the site-by-site desaturation analysis. Numbers at branches are bootstrap proportions. The scale bar indicates the number of substitutions per position.

Supplementary figs. S13-S22. Rooted maximum likelihood trees corresponding to the gene-by-gene desaturation analysis. Numbers at branches are bootstrap proportions. The scale bar indicates the number of substitutions per position.

TaxID	Species	Phylum	Class/Order
* : species retained for the construction of the final datasets (108 species)			
ARCHAEA			
1	311458 <i>Candidatus</i> Caldiarchaeum subterraneum	* Aigarchaeota	Unclassified
2	414004 <i>Cenarchaeum symbiosum</i> A	* Thaumarchaeota	Cenarchaeales
3	436308 <i>Nitrosopumilus maritimus</i> SCM1	* Thaumarchaeota	Nitrosopumilales
4	886738 <i>Candidatus</i> Nitrosoarchaeum limnia SFB1	* Thaumarchaeota	Nitrosopumilales
5	1001994 <i>Candidatus</i> Nitrosoarchaeum koreensis MY1	* Thaumarchaeota	Nitrosopumilales
6	859350 <i>Candidatus</i> Nitrosopumilus salaria BD31	* Thaumarchaeota	Nitrosopumilales
7	797209 <i>Haladaptatus paucihalophilus</i> DX253	* Euryarchaeota	Halobacteriales
8	795797 <i>Halalkalicoccus jeotgali</i> B3	Euryarchaeota	Halobacteriales
9	634497 <i>Haloarcula hispanica</i> ATCC 33960	Euryarchaeota	Halobacteriales
10	272569 <i>Haloarcula marismortui</i> ATCC 43049	Euryarchaeota	Halobacteriales
11	478009 <i>Halobacterium salinarum</i> R1	Euryarchaeota	Halobacteriales
12	309800 <i>Haloferax volcanii</i> DS2	Euryarchaeota	Halobacteriales
13	469382 <i>Halogeometricum borinquense</i> DSM 11551	* Euryarchaeota	Halobacteriales
14	485914 <i>Halomicrobium mukohataei</i> DSM 12286	Euryarchaeota	Halobacteriales
15	756883 <i>halophilic archaeon</i> DL31	Euryarchaeota	Halobacteriales
16	797210 <i>Halopiger xanaduensis</i> SH-6	Euryarchaeota	Halobacteriales
17	362976 <i>Haloquadratum walsbyi</i> DSM 16790	Euryarchaeota	Halobacteriales
18	1033806 <i>Halorhabdus tiamatea</i> SARL4B	Euryarchaeota	Halobacteriales
19	519442 <i>Halorhabdus utahensis</i> DSM 12940	* Euryarchaeota	Halobacteriales
20	416348 <i>Halorubrum lacusprofundi</i> ATCC 49239	* Euryarchaeota	Halobacteriales
21	543526 <i>Haloterrigena turkmenica</i> DSM 5511	Euryarchaeota	Halobacteriales
22	547559 <i>Natrialba magadii</i> ATCC 43099	* Euryarchaeota	Halobacteriales
23	797303 <i>Natrinema pellirubrum</i> DSM 15624	Euryarchaeota	Halobacteriales
24	797304 <i>Natronobacterium gregoryi</i> SP2	Euryarchaeota	Halobacteriales
25	348780 <i>Natronomonas pharaonis</i> DSM 2160	* Euryarchaeota	Halobacteriales
26	1072681 <i>Candidatus</i> Haloredivivus sp. G17	* Euryarchaeota	Nanohaloarchaea
27	889948 <i>Candidatus</i> Nanosalina sp. J07AB43	* Euryarchaeota	Nanohaloarchaea
28	889962 <i>Candidatus</i> Nanosalinarum sp. J07AB56	* Euryarchaeota	Nanohaloarchaea
29	410358 <i>Methanocorpusculum labreanum</i> Z	* Euryarchaeota	Methanomicrobiales
30	368407 <i>Methanoculleus marisnigri</i> JR1	* Euryarchaeota	Methanomicrobiales
31	882090 <i>Methanolinea tarda</i> NOBI-1	Euryarchaeota	Methanomicrobiales
32	937775 <i>Methanoplanus limicola</i> DSM 2279	* Euryarchaeota	Methanomicrobiales
33	679926 <i>Methanoplanus petrolearius</i> DSM 11571	Euryarchaeota	Methanomicrobiales
34	456442 <i>Methanoregula boonei</i> 6A8	Euryarchaeota	Methanomicrobiales
35	521011 <i>Methanosphaerula palustris</i> E1-9c	* Euryarchaeota	Methanomicrobiales
36	323259 <i>Methanospirillum hungatei</i> JF-1	* Euryarchaeota	Methanomicrobiales
37	259564 <i>Methanococcoides burtonii</i> DSM 6242	* Euryarchaeota	Methanosarcinales
38	644295 <i>Methanohalobium evestigatum</i> Z-7303	Euryarchaeota	Methanosarcinales
39	547558 <i>Methanohalophilus mahii</i> DSM 5219	Euryarchaeota	Methanosarcinales
40	990316 <i>Methanosaeta concilii</i> GP6	Euryarchaeota	Methanosarcinales
41	1110509 <i>Methanosaeta harundinacea</i> 6Ac	* Euryarchaeota	Methanosarcinales
42	349307 <i>Methanosaeta thermophila</i> PT	* Euryarchaeota	Methanosarcinales
43	679901 <i>Methanosalsum zhilinae</i> DSM 4017	* Euryarchaeota	Methanosarcinales
44	188937 <i>Methanosarcina acetivorans</i> C2A	Euryarchaeota	Methanosarcinales
45	269797 <i>Methanosarcina barkeri</i> str. Fusaro	Euryarchaeota	Methanosarcinales
46	192952 <i>Methanosarcina mazei</i> Go1	* Euryarchaeota	Methanosarcinales
47	351160 <i>Methanocella arvoryzae</i> MRE50	* Euryarchaeota	Methanocellales
48	1041930 <i>Methanocella conradii</i> HZ254	* Euryarchaeota	Methanocellales
49	304371 <i>Methanocella paludicola</i> SANAE	* Euryarchaeota	Methanocellales
50	224325 <i>Archaeoglobus fulgidus</i> DSM 4304	* Euryarchaeota	Archaeoglobales
51	572546 <i>Archaeoglobus profundus</i> DSM 5631	* Euryarchaeota	Archaeoglobales
52	693661 <i>Archaeoglobus veneficus</i> SNP6	* Euryarchaeota	Archaeoglobales
53	589924 <i>Ferroglobus placidus</i> DSM 10642	* Euryarchaeota	Archaeoglobales
54	439481 <i>Aciduliprofundum boonei</i> T469	* Euryarchaeota	Thermoplasmatales
55	333146 <i>Ferroplasma acidarmanus</i> fer1	* Euryarchaeota	Thermoplasmatales
56	263820 <i>Picrophilus torridus</i> DSM 9790	* Euryarchaeota	Thermoplasmatales
57	273075 <i>Thermoplasma acidophilum</i> DSM 1728	* Euryarchaeota	Thermoplasmatales
58	273116 <i>Thermoplasma volcanium</i> GSS1	* Euryarchaeota	Thermoplasmatales
59	573064 <i>Methanocaldococcus fervens</i> AG86	Euryarchaeota	Methanococcales
60	573063 <i>Methanocaldococcus infernus</i> ME	* Euryarchaeota	Methanococcales
61	243232 <i>Methanocaldococcus jannaschii</i> DSM 2661	* Euryarchaeota	Methanococcales
62	644281 <i>Methanocaldococcus</i> sp. FS406-22	Euryarchaeota	Methanococcales
63	579137 <i>Methanocaldococcus vulcanius</i> M7	Euryarchaeota	Methanococcales
64	419665 <i>Methanococcus aeolicus</i> Nankai-3	* Euryarchaeota	Methanococcales
65	402880 <i>Methanococcus maripaludis</i> C5	Euryarchaeota	Methanococcales
66	406327 <i>Methanococcus vannielii</i> SB	* Euryarchaeota	Methanococcales
67	456320 <i>Methanococcus voltae</i> A3	Euryarchaeota	Methanococcales
68	647113 <i>Methanothermococcus okinawensis</i> IH1	Euryarchaeota	Methanococcales
69	647171 <i>Methanotorris formicicus</i> Mc-S-70	Euryarchaeota	Methanococcales
70	880724 <i>Methanotorris igneus</i> Kol 5	* Euryarchaeota	Methanococcales
71	868132 <i>Methanobacterium</i> sp. AL-21	* Euryarchaeota	Methanobacteriales
72	634498 <i>Methanobrevibacter ruminantium</i> M1	Euryarchaeota	Methanobacteriales
73	420247 <i>Methanobrevibacter smithii</i> ATCC 35061	* Euryarchaeota	Methanobacteriales
74	339860 <i>Methanosphaera stadtmanae</i> DSM 3091	* Euryarchaeota	Methanobacteriales
75	79929 <i>Methanothermobacter marburgensis</i> str. Marburg	Euryarchaeota	Methanobacteriales
76	187420 <i>Methanothermobacter thermautotrophicus</i> str. Delta H	* Euryarchaeota	Methanobacteriales

77	523846	<i>Methanothermus fervidus</i> DSM 2088	* Euryarchaeota	Methanobacteriales
78	190192	<i>Methanopyrus kandleri</i> AV19	* Euryarchaeota	Methanopyrales
79	272844	<i>Pyrococcus abyssi</i> GE5	* Euryarchaeota	Thermococcales
80	186497	<i>Pyrococcus furiosus</i> DSM 3638	Euryarchaeota	Thermococcales
81	70601	<i>Pyrococcus horikoshii</i> OT3	Euryarchaeota	Thermococcales
82	342949	<i>Pyrococcus</i> sp. NA2	Euryarchaeota	Thermococcales
83	529709	<i>Pyrococcus yayanosii</i> CH1	* Euryarchaeota	Thermococcales
84	391623	<i>Thermococcus barophilus</i> MP	* Euryarchaeota	Thermococcales
85	593117	<i>Thermococcus gammatolerans</i> EJ3	* Euryarchaeota	Thermococcales
86	69014	<i>Thermococcus kodakarensis</i> KOD1	Euryarchaeota	Thermococcales
87	523849	<i>Thermococcus litoralis</i> DSM 5473	* Euryarchaeota	Thermococcales
88	523850	<i>Thermococcus onnurineus</i> NA1	Euryarchaeota	Thermococcales
89	604354	<i>Thermococcus sibiricus</i> MM 739	Euryarchaeota	Thermococcales
90	246969	<i>Thermococcus</i> sp. AM4	Euryarchaeota	Thermococcales
91	425595	<i>Candidatus</i> Micrarchaeum acidophilum ARMAN-2	* Euryarchaeota	ARMAN
92	662760	<i>Candidatus</i> Parvarchaeum acidophilum ARMAN-4	* Euryarchaeota	ARMAN
93	662762	<i>Candidatus</i> Parvarchaeum acidophilus ARMAN-5	* Euryarchaeota	ARMAN
94	228908	<i>Nanoarchaeum equitans</i> Kin4-M	* Euryarchaeota	Nanoarchaeota
95	274854	uncultured marine group II euryarchaeote	* Euryarchaeota	Unclassified
96	666510	<i>Acidilobus saccharovorans</i> 345-15	* Crenarchaeota	Desulfurococcales
97	272557	<i>Aeropyrum pernix</i> K1	* Crenarchaeota	Desulfurococcales
98	768672	<i>Desulfurococcus fermentans</i> DSM 16532	Crenarchaeota	Desulfurococcales
99	490899	<i>Desulfurococcus kamohatkinsii</i> 1221n	* Crenarchaeota	Desulfurococcales
100	765177	<i>Desulfurococcus mucosus</i> DSM 2162	Crenarchaeota	Desulfurococcales
101	415426	<i>Hyperthermus butylicus</i> DSM 5456	* Crenarchaeota	Desulfurococcales
102	453591	<i>Ignicoccus hospitalis</i> KIN4/I	* Crenarchaeota	Desulfurococcales
103	583356	<i>Ignisphaera aggregans</i> DSM 17230	* Crenarchaeota	Desulfurococcales
104	694429	<i>Pyrolobus fumarii</i> 1A	* Crenarchaeota	Desulfurococcales
105	591019	<i>Staphylothermus hellenicus</i> DSM 12710	* Crenarchaeota	Desulfurococcales
106	399550	<i>Staphylothermus marinus</i> F1	Crenarchaeota	Desulfurococcales
107	633148	<i>Thermosphaera aggregans</i> DSM 11486	* Crenarchaeota	Desulfurococcales
108	933801	<i>Acidianus hospitalis</i> W1	* Crenarchaeota	Sulfobolales
109	1006006	<i>Metallosphaera cuprina</i> Ar-4	Crenarchaeota	Sulfobolales
110	399549	<i>Metallosphaera sedula</i> DSM 5348	* Crenarchaeota	Sulfobolales
111	671065	<i>Metallosphaera yellowstonensis</i> MK1	* Crenarchaeota	Sulfobolales
112	330779	<i>Sulfolobus acidocaldarius</i> DSM 639	* Crenarchaeota	Sulfobolales
113	429572	<i>Sulfolobus islandicus</i> L.S.2.15	Crenarchaeota	Sulfobolales
114	273057	<i>Sulfolobus solfataricus</i> P2	* Crenarchaeota	Sulfobolales
115	273063	<i>Sulfolobus tokodaii</i> str. 7	* Crenarchaeota	Sulfobolales
116	397948	<i>Caldivirga maquilingensis</i> IC-167	* Crenarchaeota	Thermoproteales
117	178306	<i>Pyrobaculum aerophilum</i> str. IM2	* Crenarchaeota	Thermoproteales
118	340102	<i>Pyrobaculum arsenaticum</i> DSM 13514	* Crenarchaeota	Thermoproteales
119	410359	<i>Pyrobaculum caldifontis</i> JCM 11548	Crenarchaeota	Thermoproteales
120	384616	<i>Pyrobaculum islandicum</i> DSM 4184	* Crenarchaeota	Thermoproteales
121	698757	<i>Pyrobaculum neutrophilum</i> V24Sta	Crenarchaeota	Thermoproteales
122	1104324	<i>Pyrobaculum oguniense</i> TE7	Crenarchaeota	Thermoproteales
123	368408	<i>Pyrobaculum</i> sp. 1860	Crenarchaeota	Thermoproteales
124	444157	<i>Thermofilum pendens</i> Hrk 5	* Crenarchaeota	Thermoproteales
125	768679	<i>Thermoproteus tenax</i> Kra 1	* Crenarchaeota	Thermoproteales
126	999630	<i>Thermoproteus uzoniensis</i> 768-20	Crenarchaeota	Thermoproteales
127	572478	<i>Vulcanisaeta distributa</i> DSM 14429	* Crenarchaeota	Thermoproteales
128	985053	<i>Vulcanisaeta moutnovskia</i> 768-28	Crenarchaeota	Thermoproteales
129	374847	<i>Candidatus</i> Korarchaeum cryptofilum OPF8	* Korarchaeota	
BACTERIA				
130	134676	<i>Actinoplanes</i> sp. SE50/110	Actinobacteria	Actinobacteria
131	1133849	<i>Nocardia brasiliensis</i> ATCC 700358	* Actinobacteria	Actinobacteria
132	101510	<i>Rhodococcus jostii</i> RHA1	Actinobacteria	Actinobacteria
133	1179773	<i>Saccharothrix espanaensis</i> DSM 44229	Actinobacteria	Actinobacteria
134	463191	<i>Streptomyces sviveus</i> ATCC 29083	Actinobacteria	Actinobacteria
135	224324	<i>Aquifex aeolicus</i> VF5	Aquificae	Aquificae
136	608538	<i>Hydrogenobacter thermophilus</i> TK-6	Aquificae	Aquificae
137	123214	<i>Persephonella marina</i> EX-H1	* Aquificae	Aquificae
138	436114	<i>Sulfurihydrogenibium</i> sp. YO3AOP1	Aquificae	Aquificae
139	648996	<i>Thermovibrio ammonificans</i> HB-1	Aquificae	Aquificae
140	485918	<i>Chitinophaga pinensis</i> DSM 2588	Bacteroidetes/Chlorobi group	Bacteroidetes
141	760192	<i>Haliscomenobacter hydrossis</i> DSM 1100	Bacteroidetes/Chlorobi group	Bacteroidetes
142	700598	<i>Niastella koreensis</i> GR20-10	Bacteroidetes/Chlorobi group	Bacteroidetes
143	290317	<i>Chlorobium phaeobacteroides</i> DSM 266	* Bacteroidetes/Chlorobi group	Chlorobi
144	517418	<i>Chloroherpeton thalassium</i> ATCC 35110	Bacteroidetes/Chlorobi group	Chlorobi
145	945713	<i>Ignavibacterium album</i> JCM 16511	Bacteroidetes/Chlorobi group	Ignavibacteria
146	511051	<i>Caldisericum exile</i> AZM16c01	Caldiserica	Caldiserica
147	765952	<i>Parachlamydia acanthamoebae</i> UV-7	Chlamydiae/Verrucomicrobia group	Chlamydiae
148	331113	<i>Simkania negevensis</i> Z	Chlamydiae/Verrucomicrobia group	Chlamydiae
149	716544	<i>Waddlia chondrophila</i> WSU 86-1044	Chlamydiae/Verrucomicrobia group	Chlamydiae

150	583355	<i>Coralimargarita akajimensis</i> DSM 45221	Chlamydiae/Verrucomicrobia group	Verrucomicrobia
151	452637	<i>Opitutus terrae</i> PB90-1	* Chlamydiae/Verrucomicrobia group	Verrucomicrobia
152	926569	<i>Anaerolinea thermophila</i> UNI-1	Chloroflexi	Anaerolineae
153	926550	<i>Caldilinea aerophila</i> DSM 14535 = NBRC 104270	Chloroflexi	Caldilineae
154	326427	<i>Chloroflexus aggregans</i> DSM 9485	Chloroflexi	Chloroflexi
155	552811	<i>Dehalogenimonas lykanthroporepellens</i> BL-DC-9	* Chloroflexi	Dehalococcoidetes
156	309801	<i>Thermomicrobium roseum</i> DSM 5159	Chloroflexi	Thermomicrobia
157	653733	<i>Desulfurispirillum indicum</i> S5	* Chrysiogenetes	Chrysiogenetes
158	329726	<i>Acaryochloris marina</i> MBIC11017	Cyanobacteria	Chroococcales
159	251221	<i>Gloeobacter violaceus</i> PCC 7421	Cyanobacteria	Gloeobacteria
160	63737	<i>Nostoc punctiforme</i> PCC 73102	Cyanobacteria	Nostocales
161	179408	<i>Oscillatoria nigro-viridis</i> PCC 7112	* Cyanobacteria	Oscillatoriales
162	251229	<i>Chroococcidiopsis thermalis</i> PCC 7203	Cyanobacteria	Pleurocapsales
163	59922	<i>Prochlorococcus marinus</i> str. MIT 9303	Cyanobacteria	Prochlorales
164	768670	<i>Calditerrivibrio nitroreducens</i> DSM 19672	* Deferribacteres	Deferribacteres
165	639282	<i>Deferribacter desulfuricans</i> SSM1	Deferribacteres	Deferribacteres
166	522772	<i>Denitrovibrio acetiphilus</i> DSM 12809	Deferribacteres	Deferribacteres
167	717231	<i>Flexistipes sinusarabici</i> DSM 4947	Deferribacteres	Deferribacteres
168	937777	<i>Deinococcus peraridilitoridis</i> DSM 19664	Deinococcus-Thermus	Deinococci
169	526227	<i>Meiothermus silvanus</i> DSM 9946	* Deinococcus-Thermus	Deinococci
170	670487	<i>Oceanithermus profundus</i> DSM 14977	Deinococcus-Thermus	Deinococci
171	743525	<i>Thermus scotoductus</i> SA-01	Deinococcus-Thermus	Deinococci
172	649638	<i>Truepera radiovictrix</i> DSM 17093	Deinococcus-Thermus	Deinococci
173	309799	<i>Dictyoglomus thermophilum</i> H-6-12	Dictyoglomi	Dictyoglomina
174	515635	<i>Dictyoglomus turgidum</i> DSM 6724	* Dictyoglomi	Dictyoglomina
175	445932	<i>Elusimicrobium minutum</i> Pei191	Elusimicrobia	Elusimicrobia
176	471821	uncultured Termite group 1 bacterium phylotype Rs-D17	* Elusimicrobia	environmental samples
177	240015	<i>Acidobacterium capsulatum</i> ATCC 51196	Fibrobacteres/Acidobacteria group	Acidobacteria
178	234267	<i>Candidatus Solibacter usitatus</i> Ellin6076	* Fibrobacteres/Acidobacteria group	Acidobacteria
179	682795	<i>Granulicella mallensis</i> MP5ACTX8	Fibrobacteres/Acidobacteria group	Acidobacteria
180	926566	<i>Terriglobus roseus</i> DSM 18391	Fibrobacteres/Acidobacteria group	Acidobacteria
181	59374	<i>Fibrobacter succinogenes</i> subsp. <i>succinogenes</i> S85	Fibrobacteres/Acidobacteria group	Fibrobacteres
182	1195464	<i>Bacillus thuringiensis</i> MC28	Firmicutes	Bacilli
183	997761	<i>Paenibacillus mucilaginosus</i> K02	Firmicutes	Bacilli
184	573061	<i>Clostridium cellulovorans</i> 743B	Firmicutes	Clostridia
185	768706	<i>Desulfosporosinus orientis</i> DSM 765	* Firmicutes	Clostridia
186	650150	<i>Erysipelothrix rhusiopathiae</i>	Firmicutes	Erysipelotrichi
187	479436	<i>Veillonella parvula</i> DSM 2008	Firmicutes	Negativicutes
188	190304	<i>Fusobacterium nucleatum</i> subsp. <i>nucleatum</i> ATCC 25586	Fusobacteria	Fusobacteriia
189	572544	<i>Ilyobacter polytropus</i> DSM 2926	Fusobacteria	Fusobacteriia
190	523794	<i>Leptotrichia buccalis</i> C-1013-b	Fusobacteria	Fusobacteriia
191	526218	<i>Sebalidella termitidis</i> ATCC 33386	* Fusobacteria	Fusobacteriia
192	519441	<i>Streptobacillus moniliformis</i> DSM 12112	Fusobacteria	Fusobacteriia
193	379066	<i>Gemmatimonas aurantiaca</i> T-27	Gemmatimonadetes	Gemmatimonadetes
194	330214	<i>Candidatus Nitrospira defluvii</i>	Nitrospirae	Nitrospira
195	1048260	<i>Leptospirillum ferriphilum</i> ML-04	Nitrospirae	Nitrospira
196	1162668	<i>Leptospirillum ferrooxidans</i> C2-3	* Nitrospirae	Nitrospira
197	289376	<i>Thermodesulfobivibrio yellowstonii</i> DSM 11347	Nitrospirae	Nitrospira
198	1142394	<i>Phycisphaera mikurensis</i> NBRC 102666	Planctomycetes	Phycisphaerae
199	530564	<i>Pirellula staleyi</i> DSM 6068	Planctomycetes	Planctomycetia
200	756272	<i>Planctomyces brasiliensis</i> DSM 5305	* Planctomycetes	Planctomycetia
201	243090	<i>Rhodopirellula baltica</i> SH 1	Planctomycetes	Planctomycetia
202	886293	<i>Singulisphaera acidiphila</i> DSM 18658	Planctomycetes	Planctomycetia
203	311403	<i>Agrobacterium radiobacter</i> K84	Proteobacteria	Alphaproteobacteria
204	137722	<i>Azospirillum</i> sp. B510	* Proteobacteria	Alphaproteobacteria
205	1037409	<i>Bradyrhizobium japonicum</i> USDA 6	Proteobacteria	Alphaproteobacteria
206	460265	<i>Methylobacterium nodulans</i> ORS 2060	Proteobacteria	Alphaproteobacteria
207	216596	<i>Rhizobium leguminosarum</i> bv. <i>viciae</i> 3841	Proteobacteria	Alphaproteobacteria
208	762376	<i>Achromobacter xylosoxidans</i> A8	Proteobacteria	Betaproteobacteria
209	266265	<i>Burkholderia xenovorans</i> LB400	* Proteobacteria	Betaproteobacteria
210	1042878	<i>Cupriavidus necator</i> N-1	Proteobacteria	Betaproteobacteria
211	398578	<i>Delftia acidovorans</i> SPH-1	Proteobacteria	Betaproteobacteria
212	381666	<i>Ralstonia eutropha</i> H16	Proteobacteria	Betaproteobacteria
213	572480	<i>Arcobacter nitrofigilis</i> DSM 7299	* Proteobacteria	delta/epsilon subdivisions
214	706587	<i>Desulfomonile tiedjei</i> DSM 6799	* Proteobacteria	delta/epsilon subdivisions
215	502025	<i>Haliangium ochraceum</i> DSM 14365	Proteobacteria	delta/epsilon subdivisions
216	246197	<i>Myxococcus xanthus</i> DK 1622	Proteobacteria	delta/epsilon subdivisions
217	709032	<i>Sulfuricurvum kujijense</i> DSM 16994	Proteobacteria	delta/epsilon subdivisions
218	349521	<i>Hahella chejuensis</i> KCTC 2396	* Proteobacteria	Gammaproteobacteria
219	1191061	<i>Klebsiella oxytoca</i> E718	Proteobacteria	Gammaproteobacteria
220	592316	<i>Pantoea</i> sp. <i>At-9b</i>	Proteobacteria	Gammaproteobacteria
221	220664	<i>Pseudomonas protegens</i> Pf-5	Proteobacteria	Gammaproteobacteria

222	338187	<i>Vibrio harveyi</i> ATCC BAA-1116	Proteobacteria	Gammaproteobacteria
223	573825	<i>Leptospira interrogans</i> serovar Lai str. IPAV	* Spirochaetes	Spirochaetia
224	158190	<i>Sphaerochaeta pleomorpha</i> str. Grapes	Spirochaetes	Spirochaetia
225	573413	<i>Spirochaeta smaragdinae</i> DSM 11293	Spirochaetes	Spirochaetia
226	545694	<i>Treponema primitia</i> ZAS-2	Spirochaetes	Spirochaetia
227	869212	<i>Turneriella parva</i> DSM 21527	Spirochaetes	Spirochaetia
228	572547	<i>Aminobacterium colombiense</i> DSM 12261	Synergistetes	Synergistia
229	584708	<i>Aminomonas paucivorans</i> DSM 12260	Synergistetes	Synergistia
230	891968	<i>Anaerobaculum mobile</i> DSM 13181	* Synergistetes	Synergistia
231	580340	<i>Thermovirga lienii</i> DSM 17291	Synergistetes	Synergistia
232	441768	<i>Acholeplasma laidlawii</i> PG-8A	Tenericutes	Mollicutes
233	322098	<i>Aster yellows witches'-broom phytoplasma</i> AYWB	Tenericutes	Mollicutes
234	265311	<i>Mesoplasma florum</i> L1	Tenericutes	Mollicutes
235	272633	<i>Mycoplasma penetrans</i> HF-2	Tenericutes	Mollicutes
236	565575	<i>Ureaplasma urealyticum</i> serovar 10 str. ATCC 33699	Tenericutes	Mollicutes
237	667014	<i>Thermodesulfatator indicus</i> DSM 15286	* Thermodesulfobacteria	Thermodesulfobacteria
238	795359	<i>Thermodesulfobacterium</i> sp. OPB45	Thermodesulfobacteria	Thermodesulfobacteria
239	521045	<i>Kosmotoga olearia</i> TBF 19.5.1	Thermotogae	Thermotogae
240	443254	<i>Marinitoga piezophila</i> KA3	Thermotogae	Thermotogae
241	660470	<i>Mesotoga prima</i> MesG1.Ag.4.2	Thermotogae	Thermotogae
242	403833	<i>Petrotoga mobilis</i> SJ95	Thermotogae	Thermotogae
243	126740	<i>Thermotoga</i> sp. RQ2	* Thermotogae	Thermotogae
244	880073	<i>Caldithrix abyssi</i> DSM 13497	unclassified Bacteria	Caldithrix
245	671143	<i>Candidatus Methyloirabilis oxyfera</i>	* unclassified Bacteria	candidate division NC10
246	525904	<i>Thermobaculum terrenum</i> ATCC BAA-798	* unclassified Bacteria	Thermobaculum

Annotation	Reference gi	Reference accession	COG number	COG class	COG definition
argininosuccinate synthase	161528794	YP_001582620.1	COG0137	E	Amino acid transport and metabolism
homoserine kinase	161527553	YP_001581379.1	COG0083	E	Amino acid transport and metabolism
homoserine dehydrogenase	161528204	YP_001582030.1	COG0460	E	Amino acid transport and metabolism
aspartate kinase	161529262	YP_001583088.1	COG0527	E	Amino acid transport and metabolism
aspartate carbamoyltransferase	161529194	YP_001583020.1	COG0540	F	Nucleotide transport and metabolism
phosphoribosylformylglycinamide cyclo-ligase	161528090	YP_001581916.1	COG0150	F	Nucleotide transport and metabolism
uridylyate kinase putative	161529213	YP_001583039.1	COG0528	F	Nucleotide transport and metabolism
phosphoribosylformylglycinamide synthase II	161528288	YP_001582114.1	COG0046	F	Nucleotide transport and metabolism
adenylosuccinate lyase	161528963	YP_001582789.1	COG0015	F	Nucleotide transport and metabolism
RdgB/HAM1 family non-canonical purine NTP pyrophosphatase	161529038	YP_001582864.1	COG0127	F	Nucleotide transport and metabolism
glutamine amidotransferase class-II	161528289	YP_001582115.1	COG0034	F	Nucleotide transport and metabolism
phosphoglycerate kinase	161528008	YP_001581834.1	COG0126	G	Carbohydrate transport and metabolism
phosphopantothenoylcysteine decarboxylase/phosphopantothenate--cysteine ligase	161529229	YP_001583055.1	COG0452	H	Coenzyme transport and metabolism
pyridoxine biosynthesis protein	161528987	YP_001582813.1	COG0214	H	Coenzyme transport and metabolism
UbiD family decarboxylase	161529105	YP_001582931.1	COG0043	H	Coenzyme transport and metabolism
porphobilinogen deaminase	118576540	YP_876283.1	COG0181	H	Coenzyme transport and metabolism
glutamate-1-semialdehyde-2 1-aminomutase	161527998	YP_001581824.1	COG0001	H	Coenzyme transport and metabolism
molybdenum cofactor biosynthesis protein C	315427188	BAJ48802.1	COG0315	H	Coenzyme transport and metabolism
conserved hypothetical protein	315426923	BAJ48542.1	COG3425	I	Lipid transport and metabolism
beta-lactamase domain-containing protein	161529047	YP_001582873.1	COG1236	J	Translation, ribosomal structure and biogenesis
glutamyl-tRNA(Gln) amidotransferase subunit E	161527611	YP_001581437.1	COG2511	J	Translation, ribosomal structure and biogenesis
tRNA-guanine transglycosylase	161529044	YP_001582870.1	COG0343	J	Translation, ribosomal structure and biogenesis
glutamyl-tRNA(Gln) amidotransferase B subunit	161528376	YP_001582202.1	COG0064	J	Translation, ribosomal structure and biogenesis
exosome complex exonuclease 1	161527940	YP_001581766.1	COG0689	J	Translation, ribosomal structure and biogenesis
MiaB-like tRNA modifying enzyme	161528769	YP_001582595.1	COG0621	J	Translation, ribosomal structure and biogenesis
phenylalanyl-tRNA synthetase alpha subunit	161528997	YP_001582823.1	COG0016	J	Translation, ribosomal structure and biogenesis
prolyl-tRNA synthetase	315425642	BAJ47301.1	COG0442	J	Translation, ribosomal structure and biogenesis
histidyl-tRNA synthetase	315425704	BAJ47360.1	COG0124	J	Translation, ribosomal structure and biogenesis
5-nucleotidase SurE	315426747	BAJ48371.1	COG0522	J	Translation, ribosomal structure and biogenesis
peptidase M50	161527799	YP_001581625.1	COG0750	M	Cell wall/membrane/envelope biogenesis
cytidyltransferase-like protein	161528806	YP_001582632.1	COG0615	MI	Cell wall/membrane/envelope biogenesis
metalloendopeptidase glycoprotease family	161529041	YP_001582867.1	COG0533	O	Post-translational modification, protein turnover, and chaperones
hydrogenase maturation protein HypF	315426449	BAJ48087.1	COG0068	O	Post-translational modification, protein turnover, and chaperones
hydrogenase expression formation protein HypD	315426462	BAJ48095.1	COG0409	O	Post-translational modification, protein turnover, and chaperones
GTP1/OBG protein	161529113	YP_001582939.1	COG2262	R	General function prediction only
LPPG:FO 2-phospho-L-lactate transferase	161528134	YP_001581960.1	COG0391	S	Function unknown
GTP-binding signal recognition particle	161528039	YP_001581865.1	COG0541	U	Intracellular trafficking, secretion, and vesicular transport
signal recognition particle receptor	315425933	BAJ47583.1	COG0552	U	Intracellular trafficking, secretion, and vesicular transport
50S ribosomal protein L1	161527890	YP_001581716			
50S ribosomal protein L2	161527614	YP_001581440			
50S ribosomal protein L3P	161528317	YP_001582143			
50S ribosomal protein L4P	161528316	YP_001582142			
50S ribosomal protein L5	161528305	YP_001582131			
50S ribosomal protein L6	161528302	YP_001582128			
acidic ribosomal protein P0 - L10	161527889	YP_001581715			
50S ribosomal protein L10e	161527957	YP_001581783			
50S ribosomal protein L11	161527893	YP_001581719			
50S ribosomal protein L13	161527934	YP_001581760			
50S ribosomal protein L14	161528308	YP_001582134			
50S ribosomal protein L15	161527909	YP_001581735			
50S ribosomal protein L18	161527906	YP_001581732			
50S ribosomal protein L22	161528313	YP_001582139			
50S ribosomal protein L25	161528315	YP_001582141			

KOW domain-containing protein – L24	161528307	YP_001582133
50S ribosomal protein L29	161528311	YP_001582137
50S ribosomal protein L30	161527908	YP_001581734
30S ribosomal protein S2	161527824	YP_001581650
30S ribosomal protein S3	161528312	YP_001582138
30S ribosomal protein S4	161527832	YP_001581658
30S ribosomal protein S5	161527907	YP_001581733
30S ribosomal protein S7	161527863	YP_001581689
30S ribosomal protein S8	161528303	YP_001582129
30S ribosomal protein S9	161527933	YP_001581759
30S ribosomal protein S10	161528541	YP_001582367
30S ribosomal protein S11	161528958	YP_001582784
30S ribosomal protein S12	161527862	YP_001581688
30S ribosomal protein S13	161527833	YP_001581659
30S ribosomal protein S15	161529016	YP_001582842
30S ribosomal protein S17	161528309	YP_001582135
30S ribosomal protein S19	161528314	YP_001582140

Nitrosopumilus maritimus SCM1 versus *Methanoculleus marisnigri* JR1
Nitrosopumilus maritimus SCM1 versus *Thermococcus gammatolerans* EJ3
Nitrosopumilus maritimus SCM1 versus *Metallosphaera sedula* DSM 5348
Nitrosopumilus maritimus SCM1 versus *Pyrobaculum islandicum* DSM 4184
Methanoculleus marisnigri JR1 versus *Thermococcus gammatolerans* EJ3
Methanoculleus marisnigri JR1 versus *Metallosphaera sedula* DSM 5348
Methanoculleus marisnigri JR1 versus *Pyrobaculum islandicum* DSM 4184
Thermococcus gammatolerans EJ3 versus *Metallosphaera sedula* DSM 5348
Thermococcus gammatolerans EJ3 versus *Pyrobaculum islandicum* DSM 4184
Metallosphaera sedula DSM 5348 versus *Pyrobaculum islandicum* DSM 4184
Nitrosopumilus maritimus SCM1 versus *Persephonella marina* EX-H1
Nitrosopumilus maritimus SCM1 versus *Azospirillum* sp. B510
Nitrosopumilus maritimus SCM1 versus *Oscillatoria nigro-viridis* PCC 7112
Methanoculleus marisnigri JR1 versus *Persephonella marina* EX-H1
Methanoculleus marisnigri JR1 versus *Azospirillum* sp. B510
Methanoculleus marisnigri JR1 versus *Oscillatoria nigro-viridis* PCC 7112
Metallosphaera sedula DSM 5348 versus *Persephonella marina* EX-H1
Metallosphaera sedula DSM 5348 versus *Azospirillum* sp. B510
Metallosphaera sedula DSM 5348 versus *Oscillatoria nigro-viridis* PCC 7112
Persephonella marina EX-H1 versus *Azospirillum* sp. B510
Persephonella marina EX-H1 versus *Oscillatoria nigro-viridis* PCC 7112
Azospirillum sp. B510 versus *Oscillatoria nigro-viridis* PCC 7112



Figure S1
108 species – 9450 positions

0.3

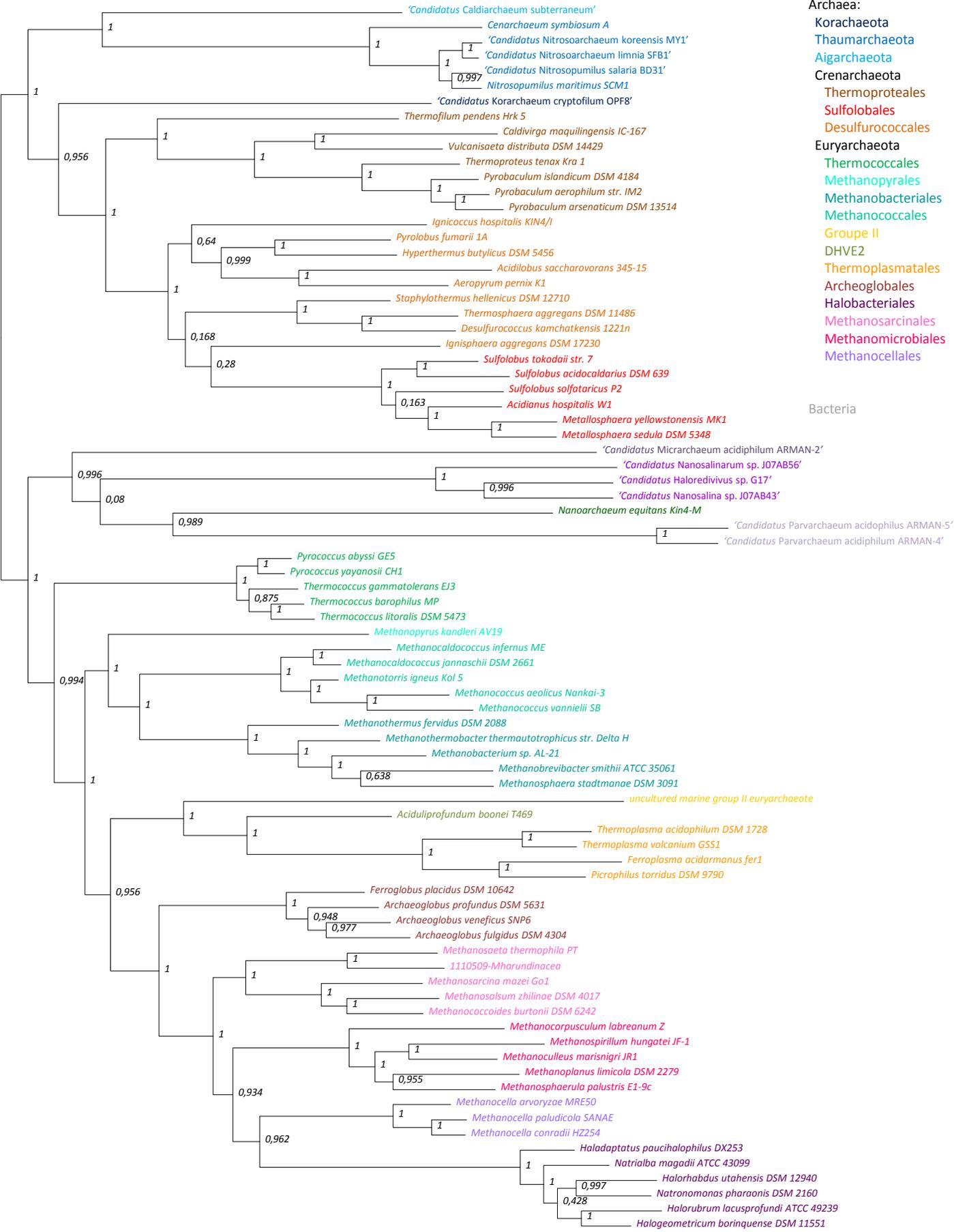
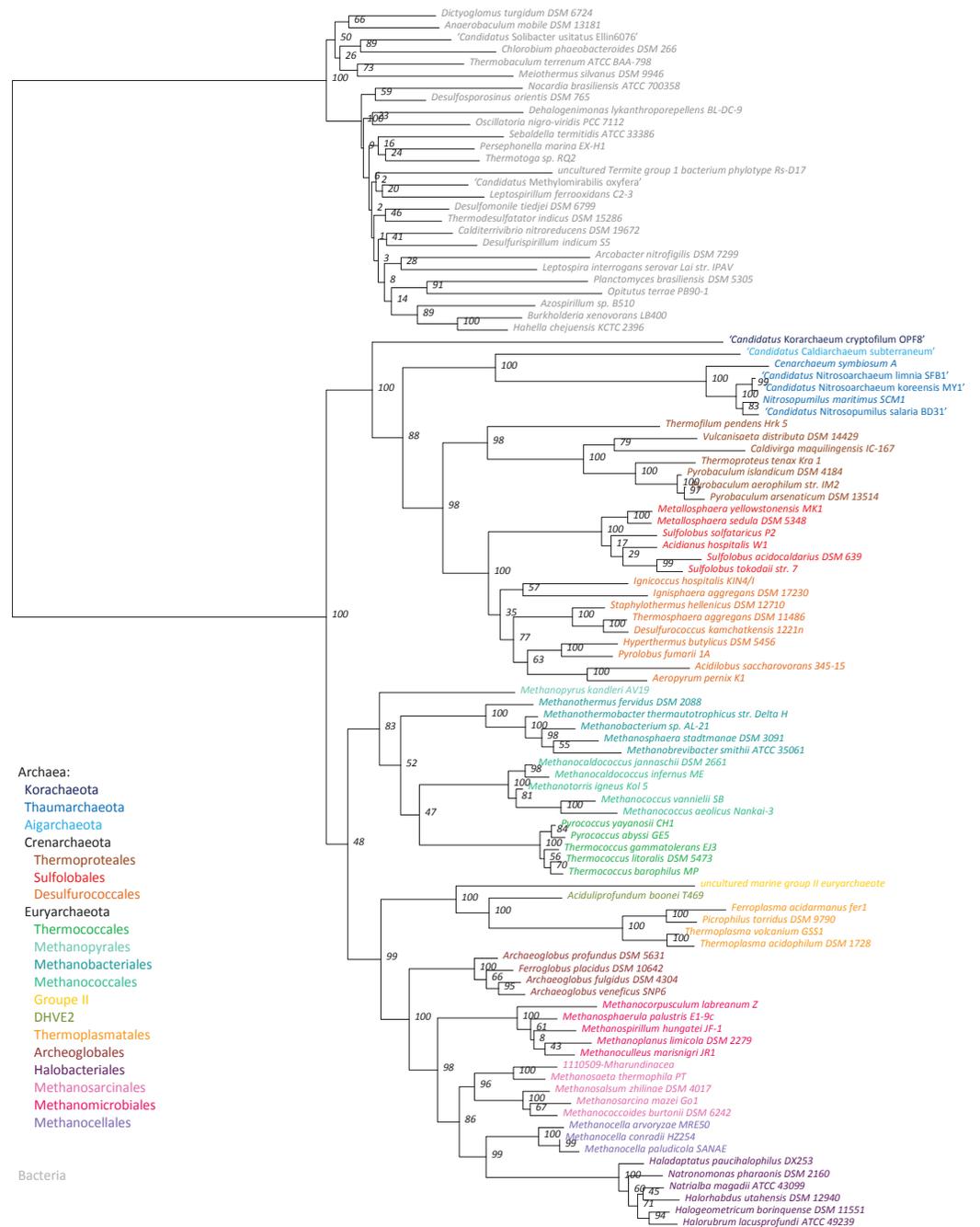


Figure S2
81 species – 9450 positions

0.2



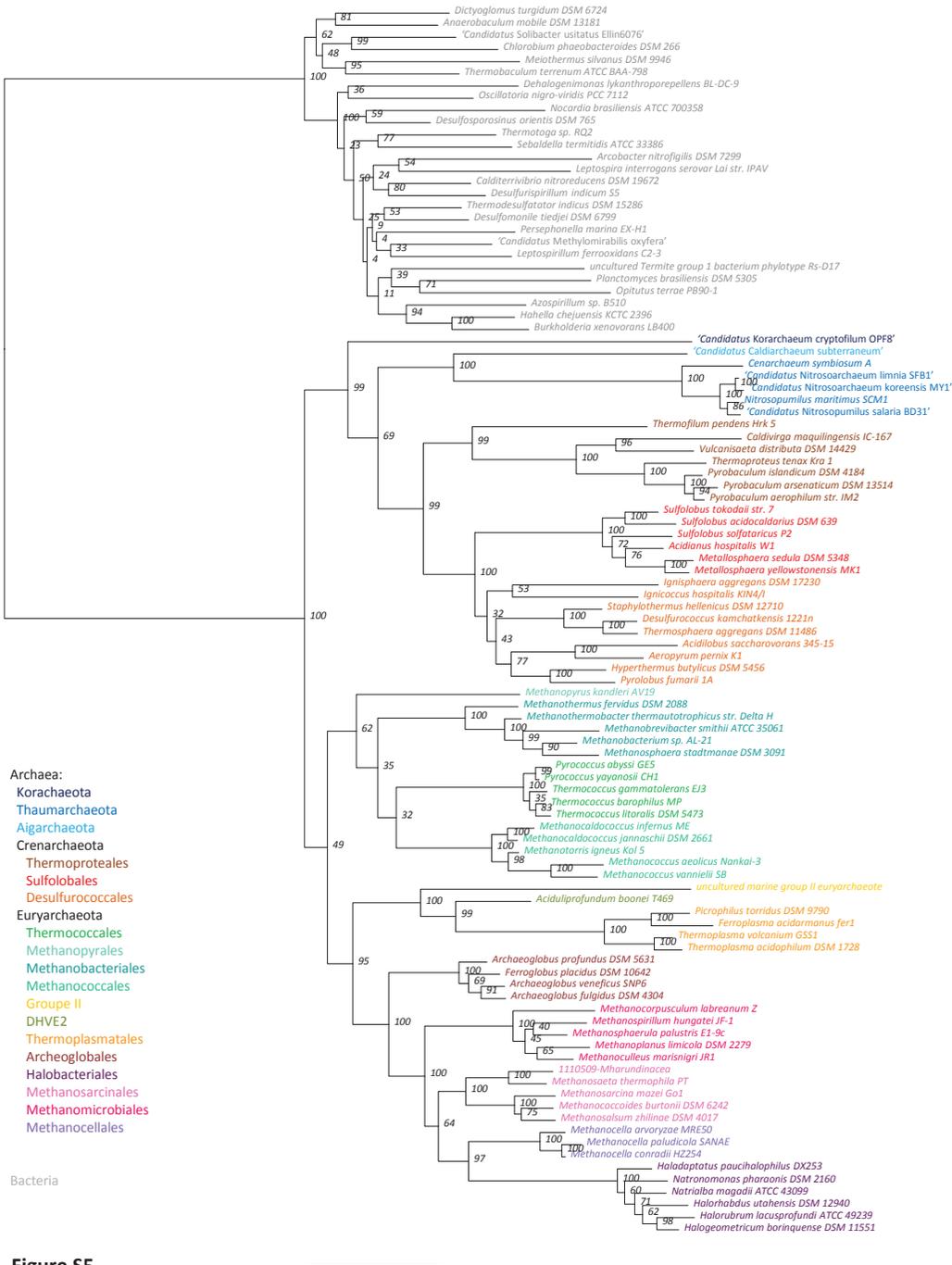


Figure S5
101 species - 2905 positions

0.07



Figure S6
101 species - 3803 positions

0.07

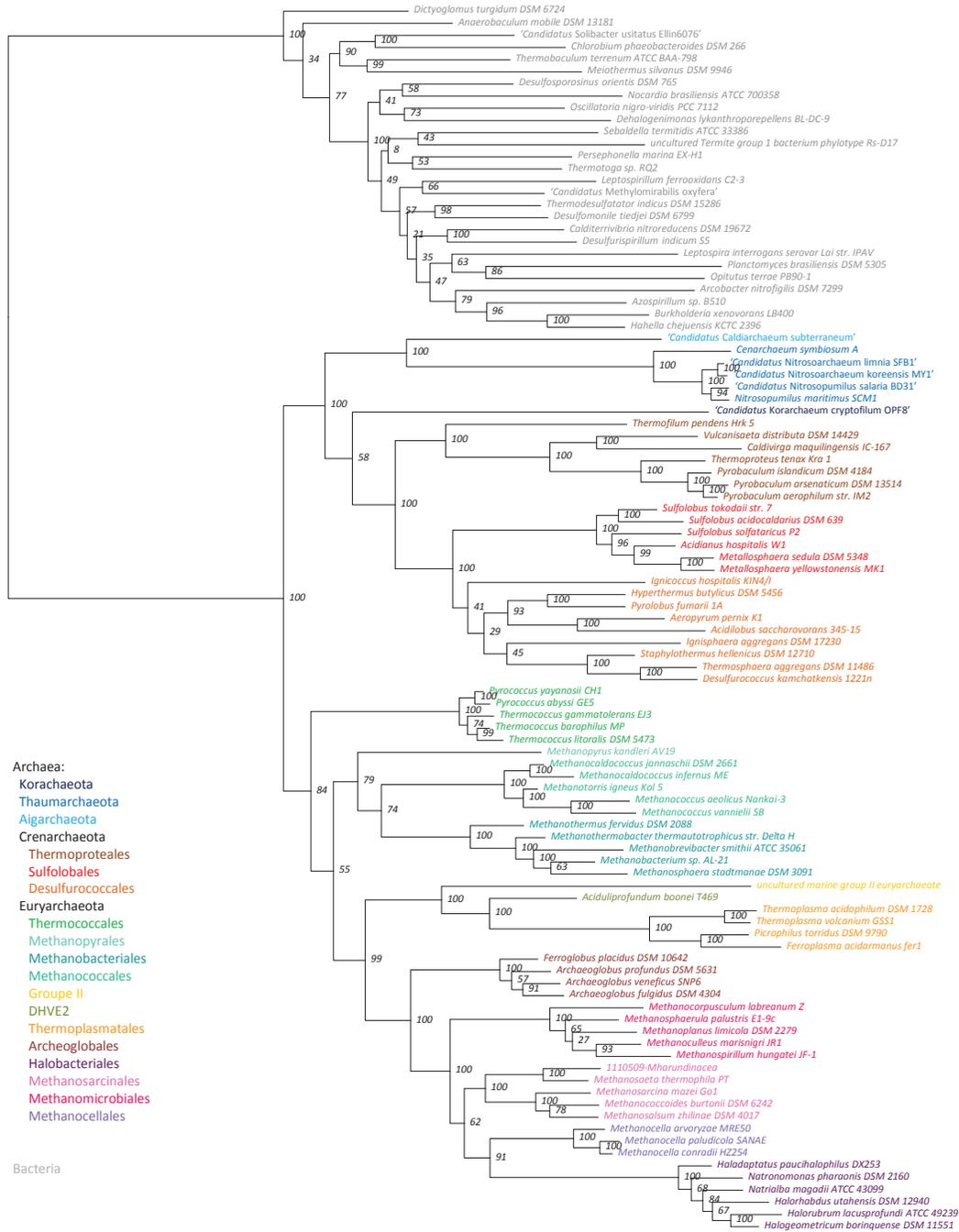


Figure S7
101 species - 4727 positions

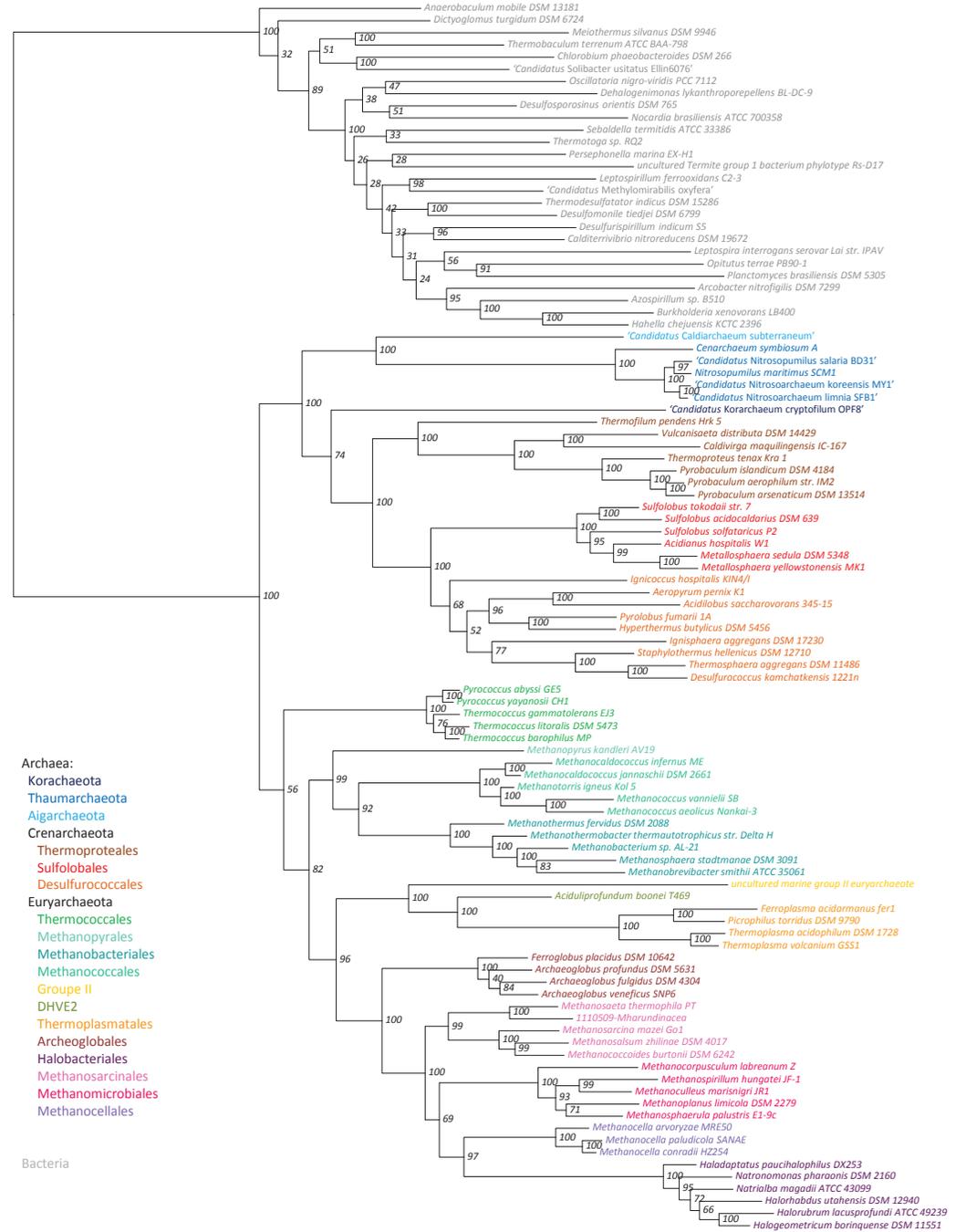


Figure S8
101 species - 5860 positions

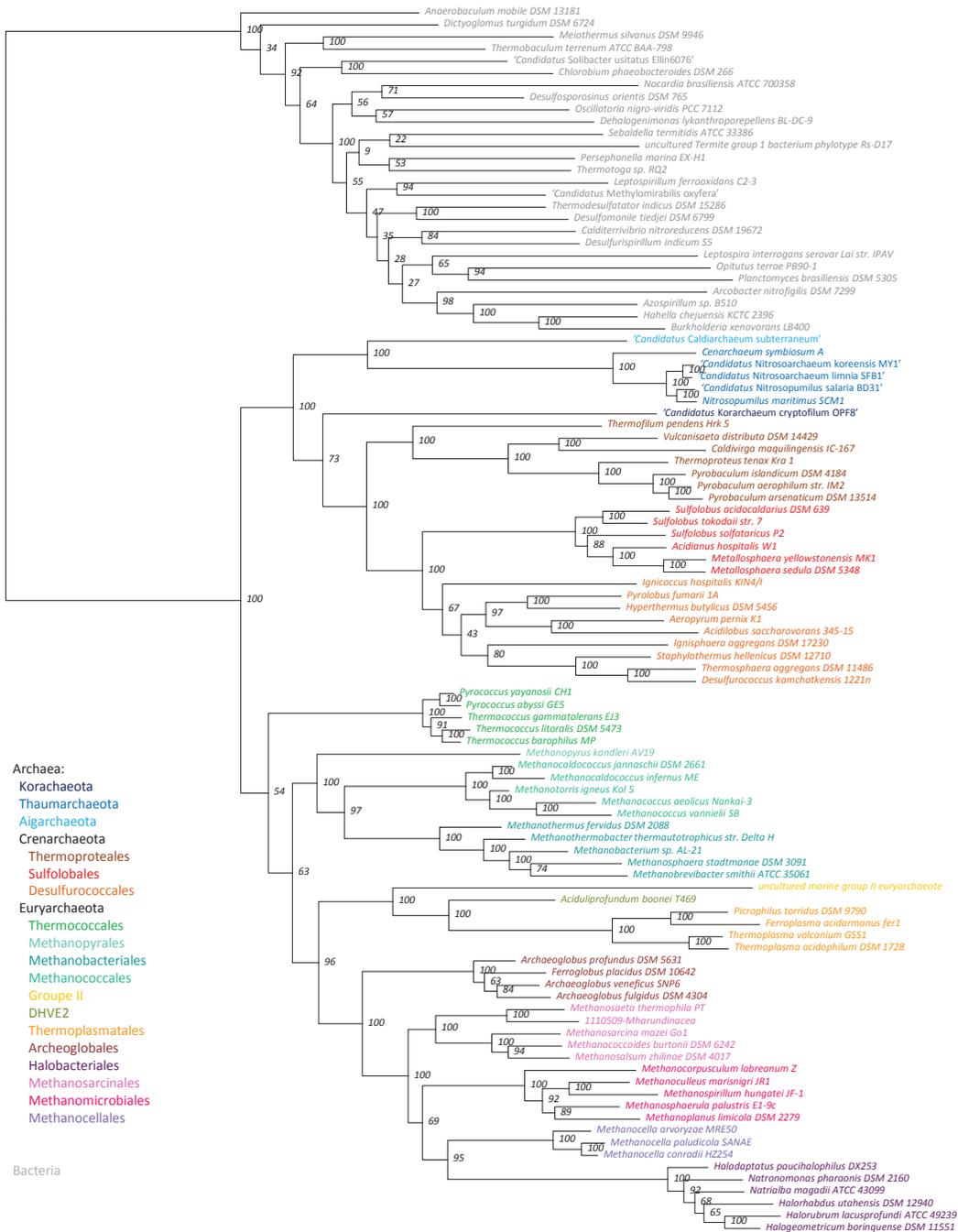


Figure S9
 101 species - 6939 positions

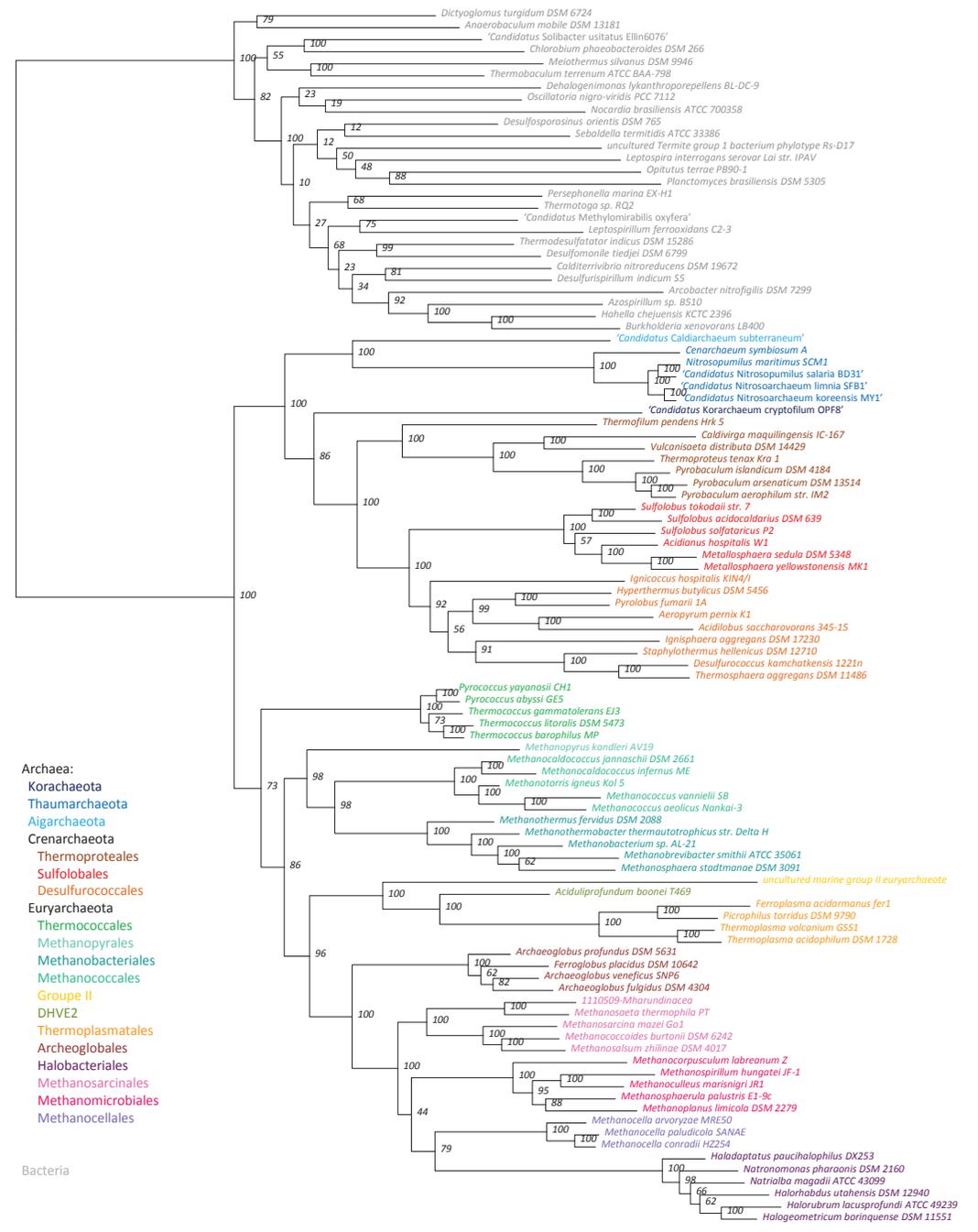


Figure S10
 101 species - 7978 positions

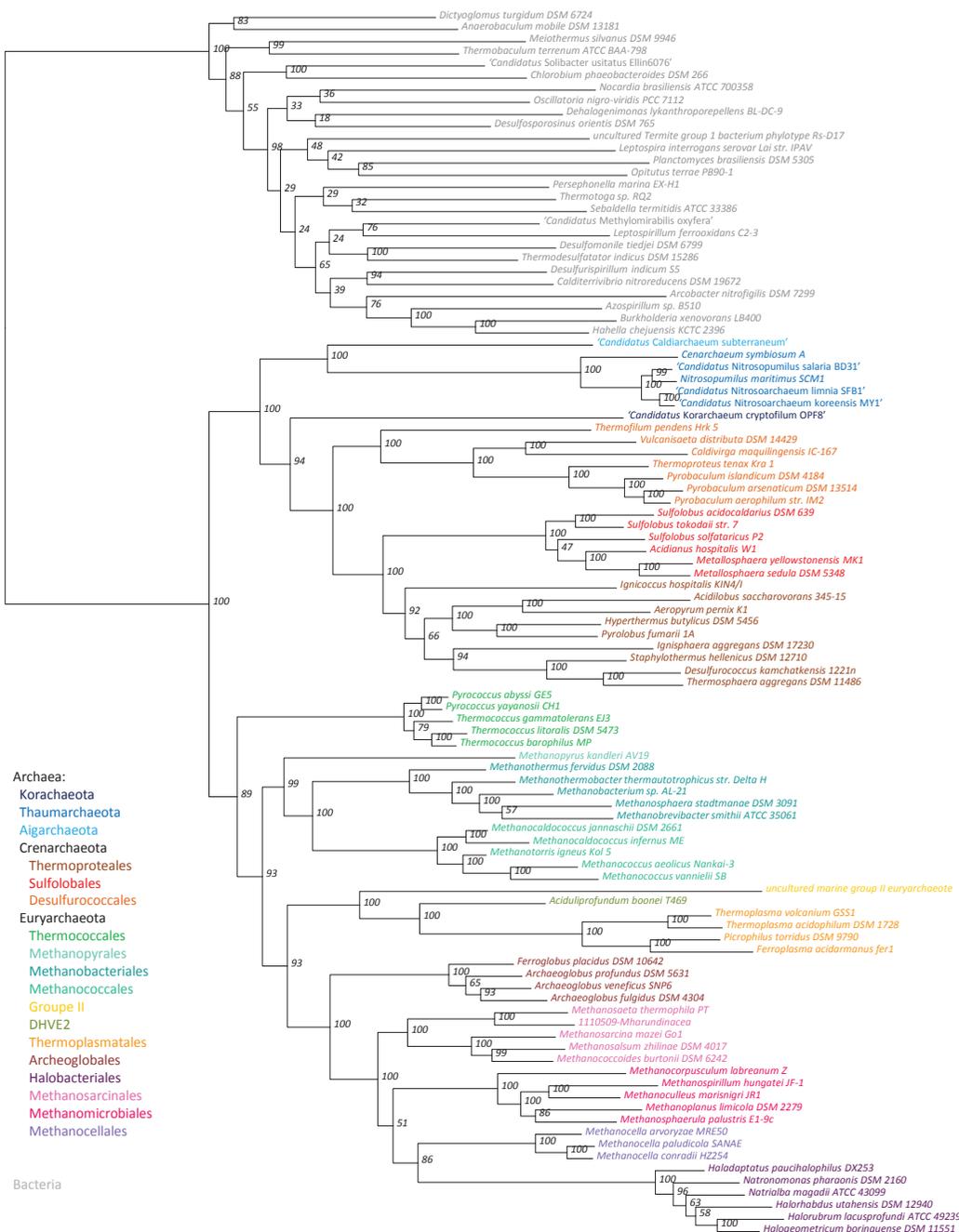


Figure S11
101 species - 8980 positions

0.2

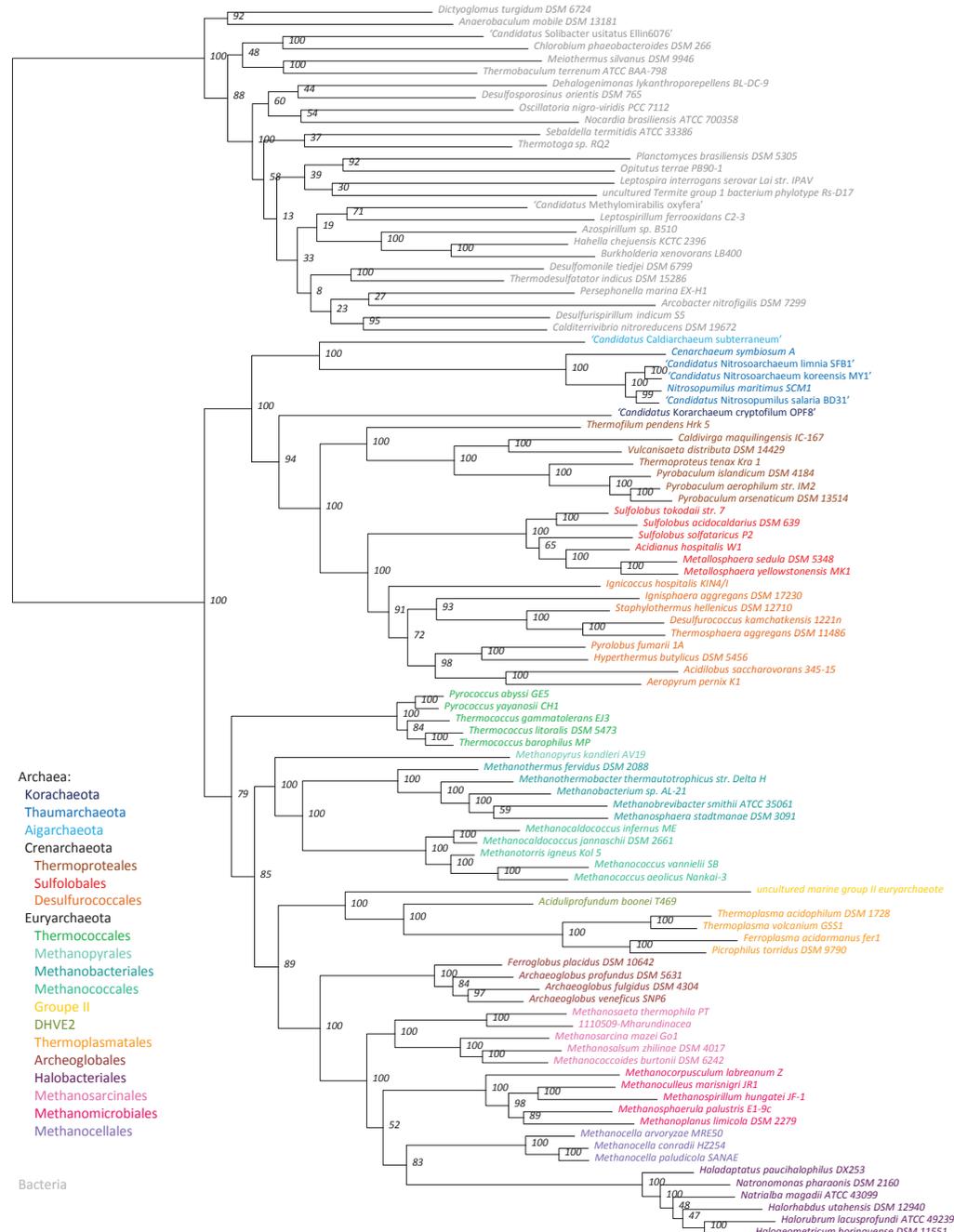


Figure S12
101 species - 9450 positions

0.3

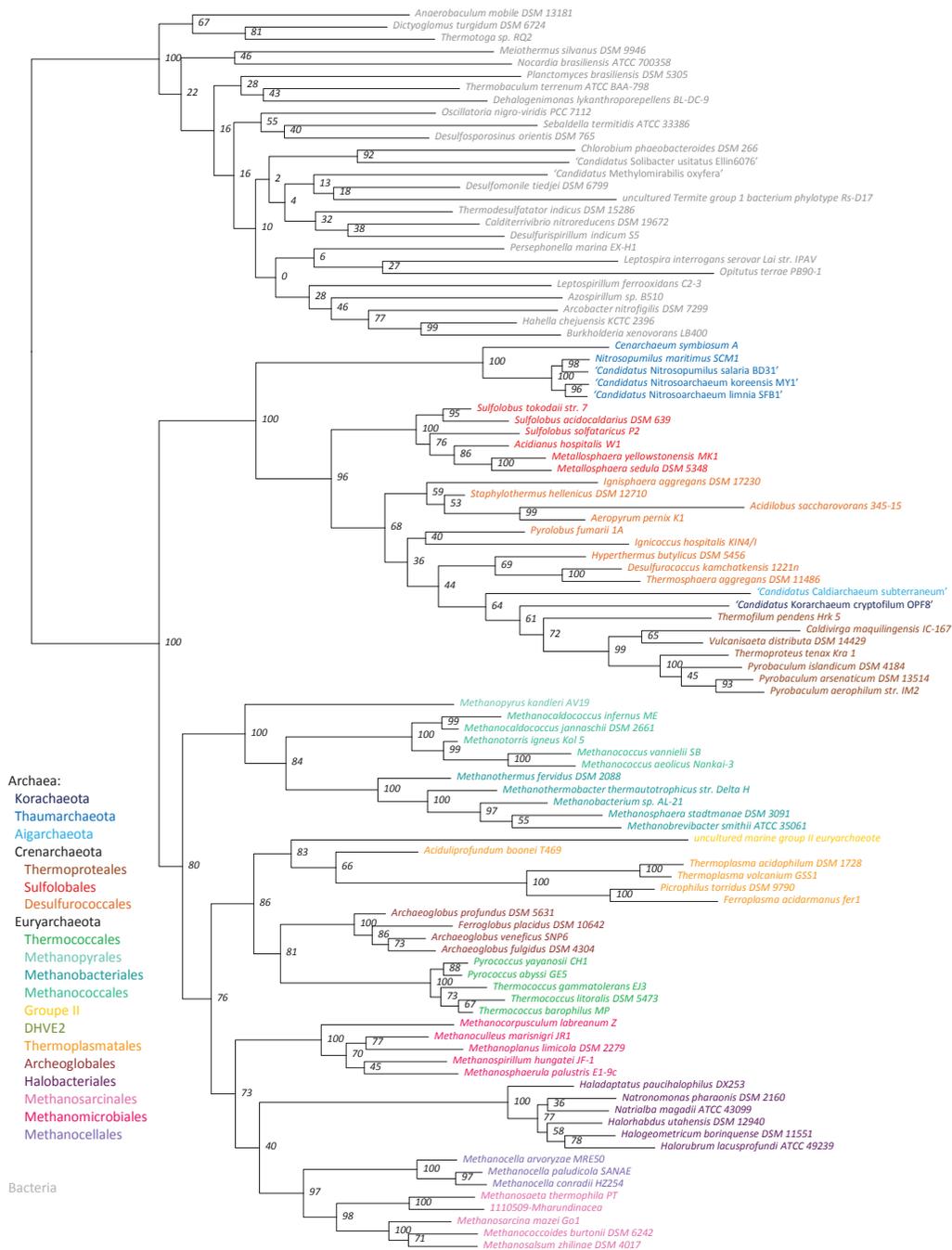


Figure S13
101 species - 870 positions

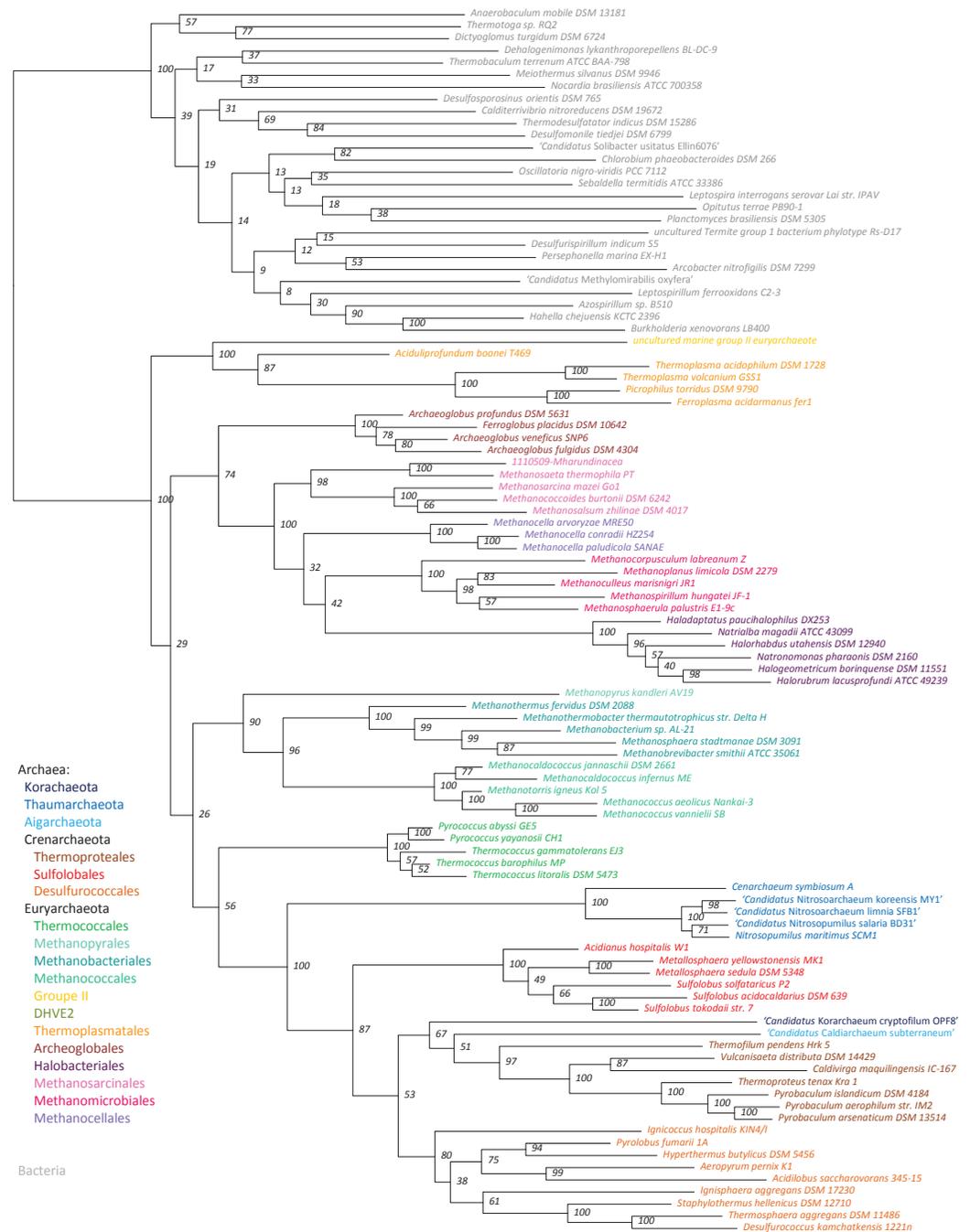


Figure S14
101 species - 1926 positions

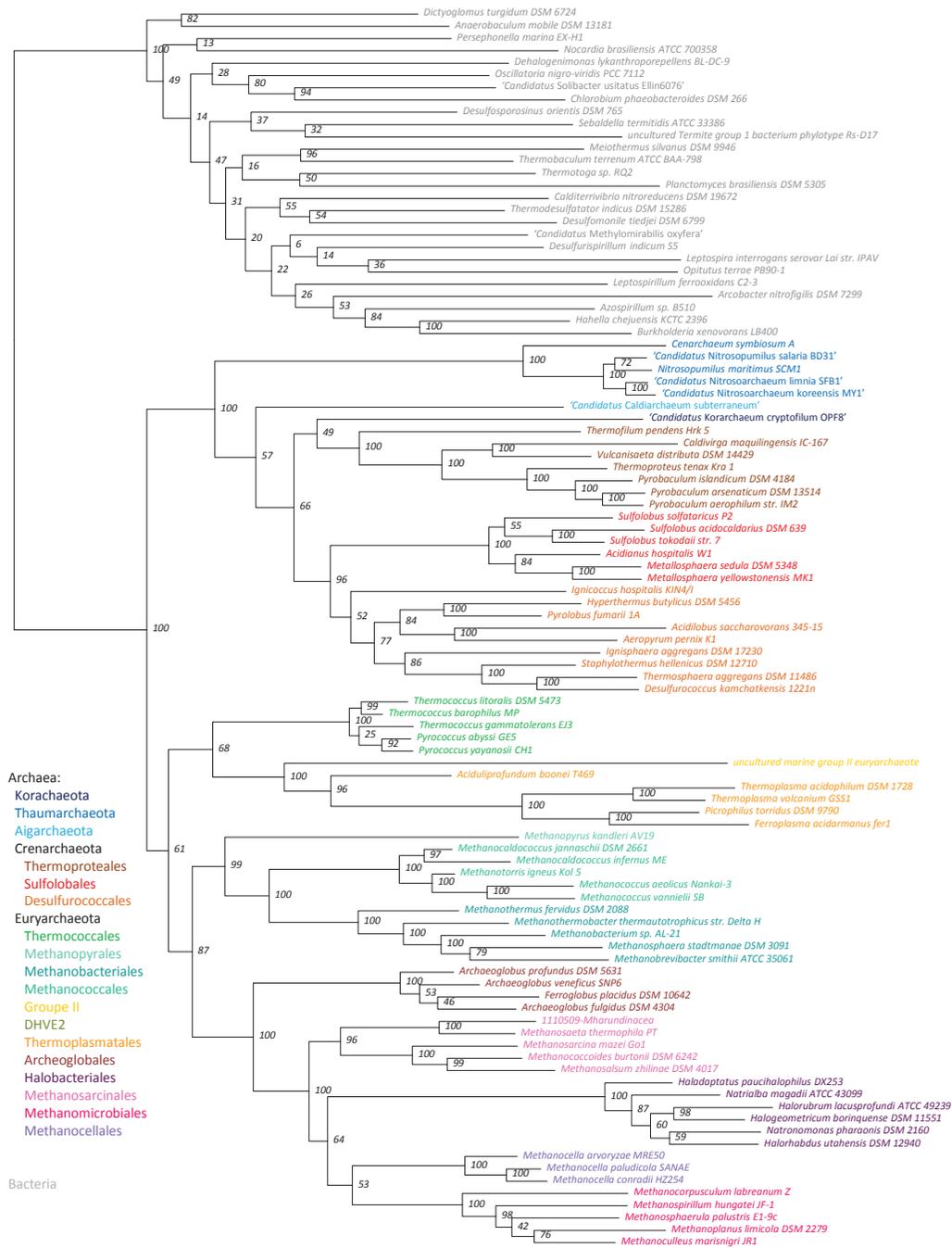


Figure S15
 101 species - 3172 positions

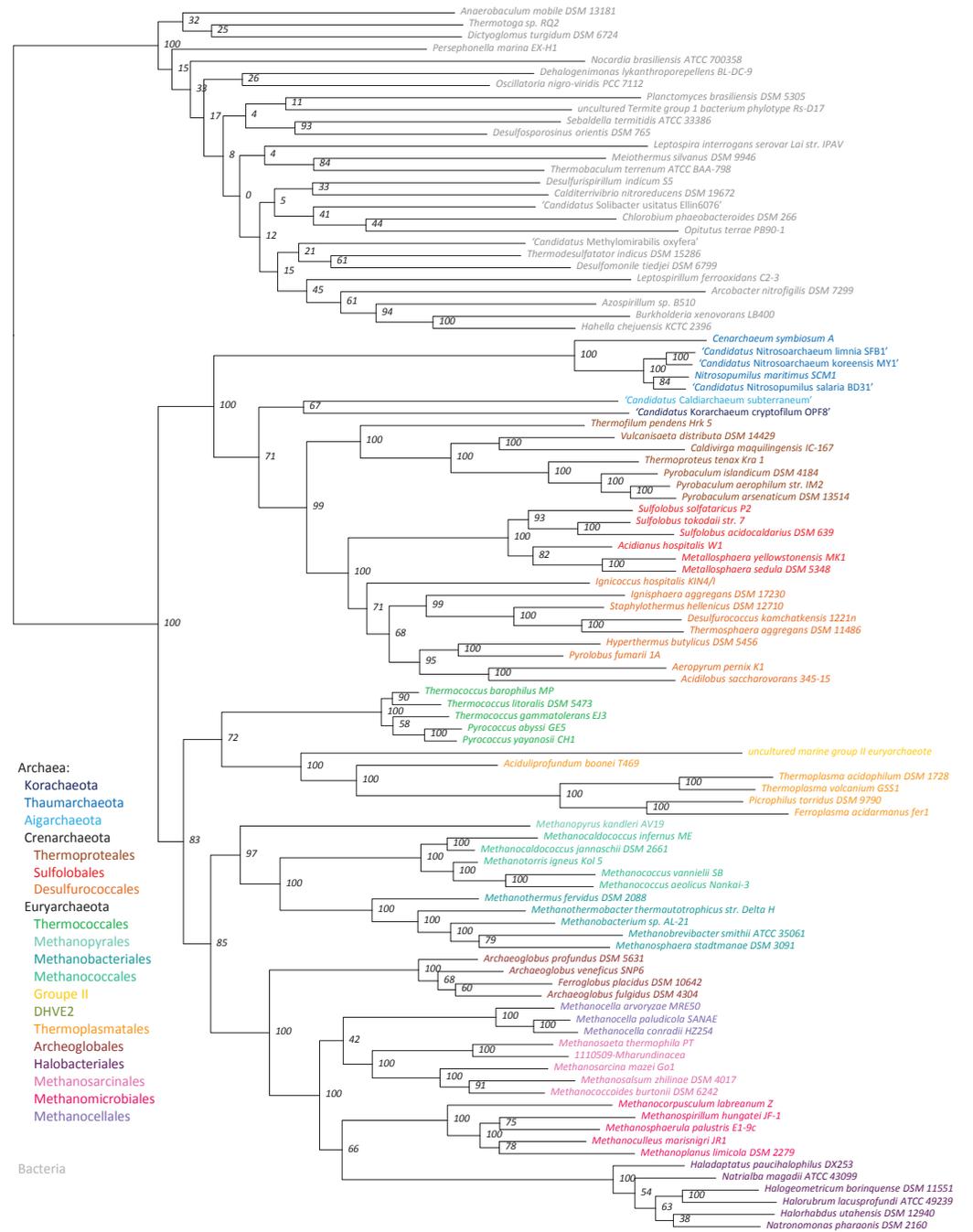


Figure S16
 101 species - 4175 positions

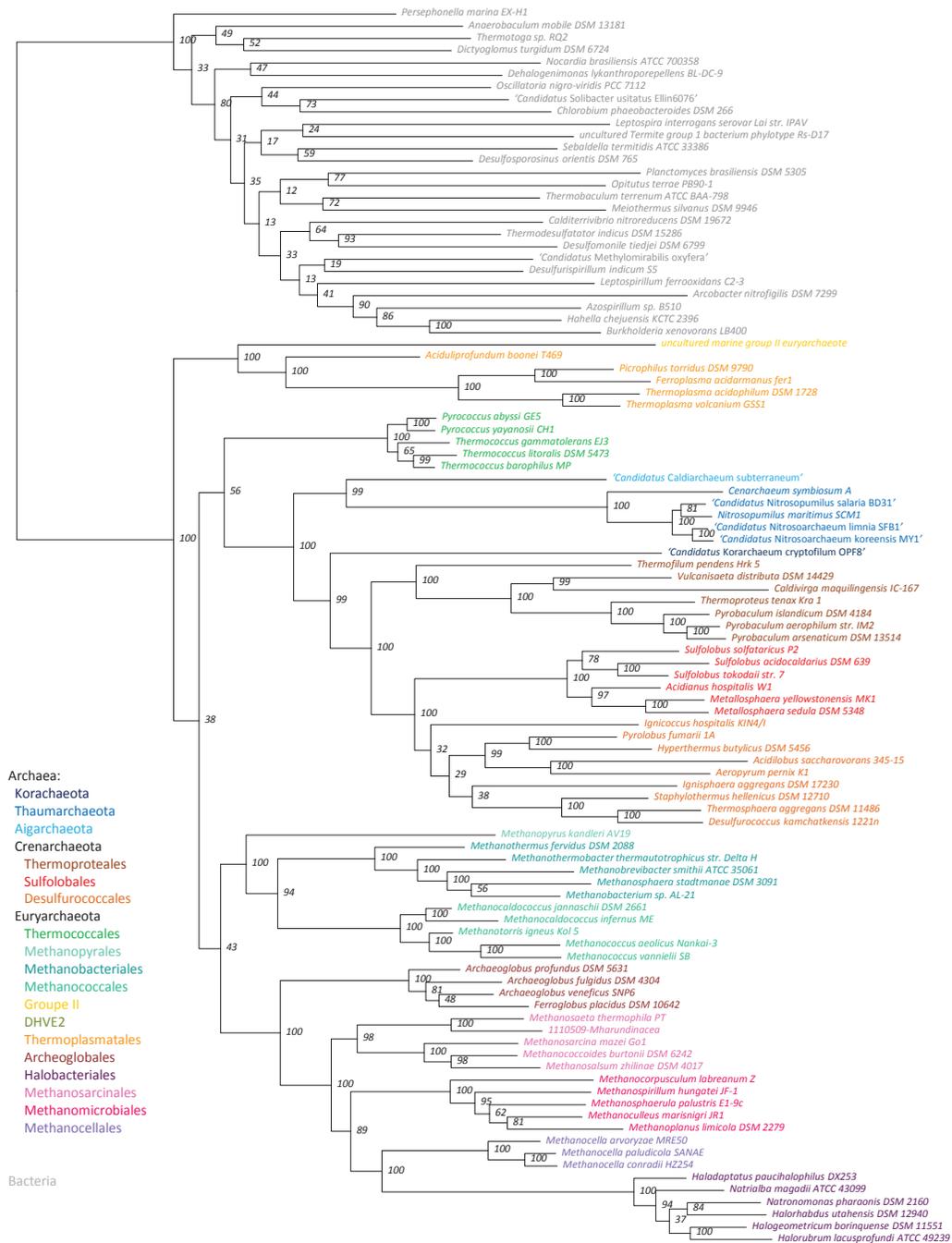


Figure S17
101 species - 4902 positions

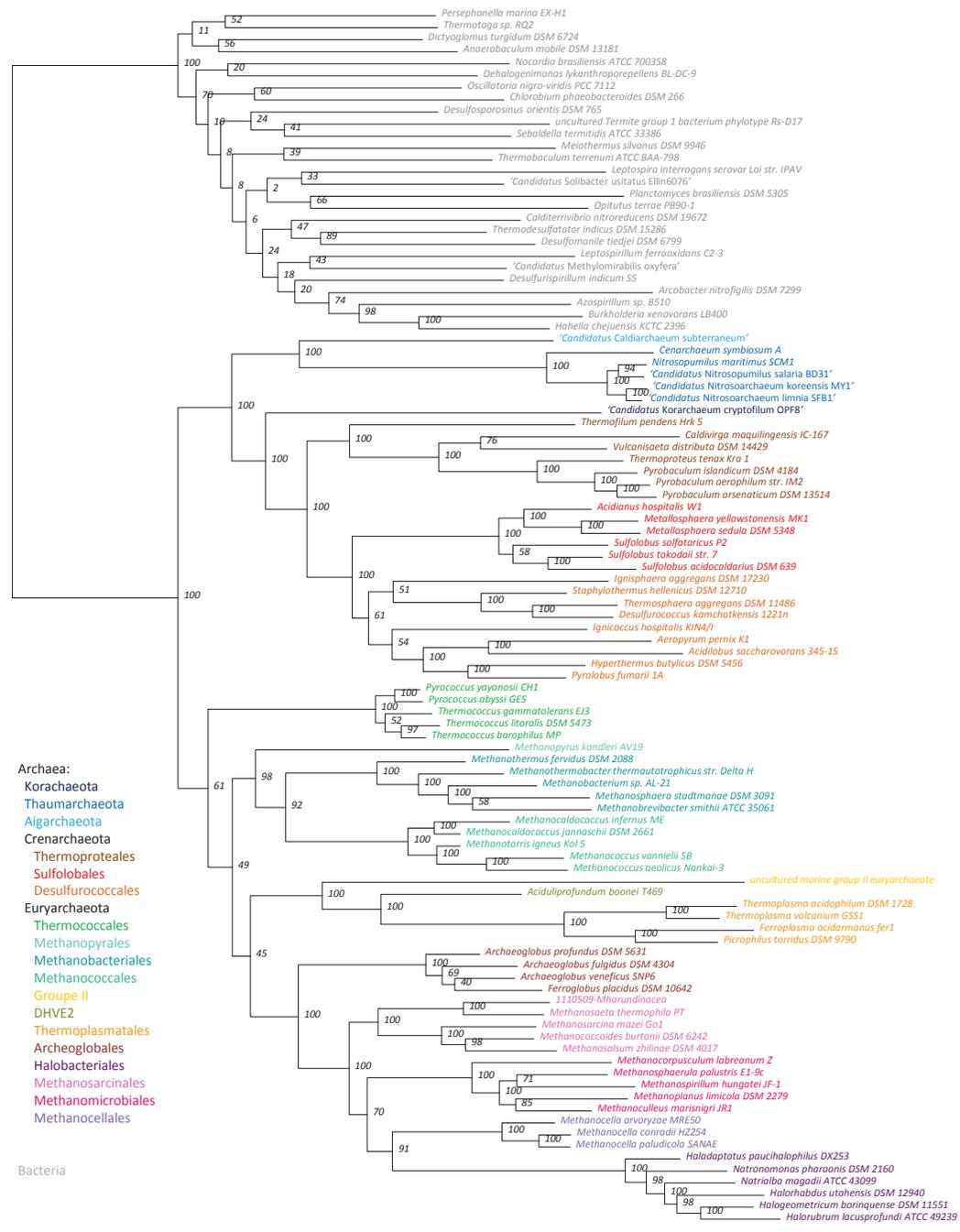
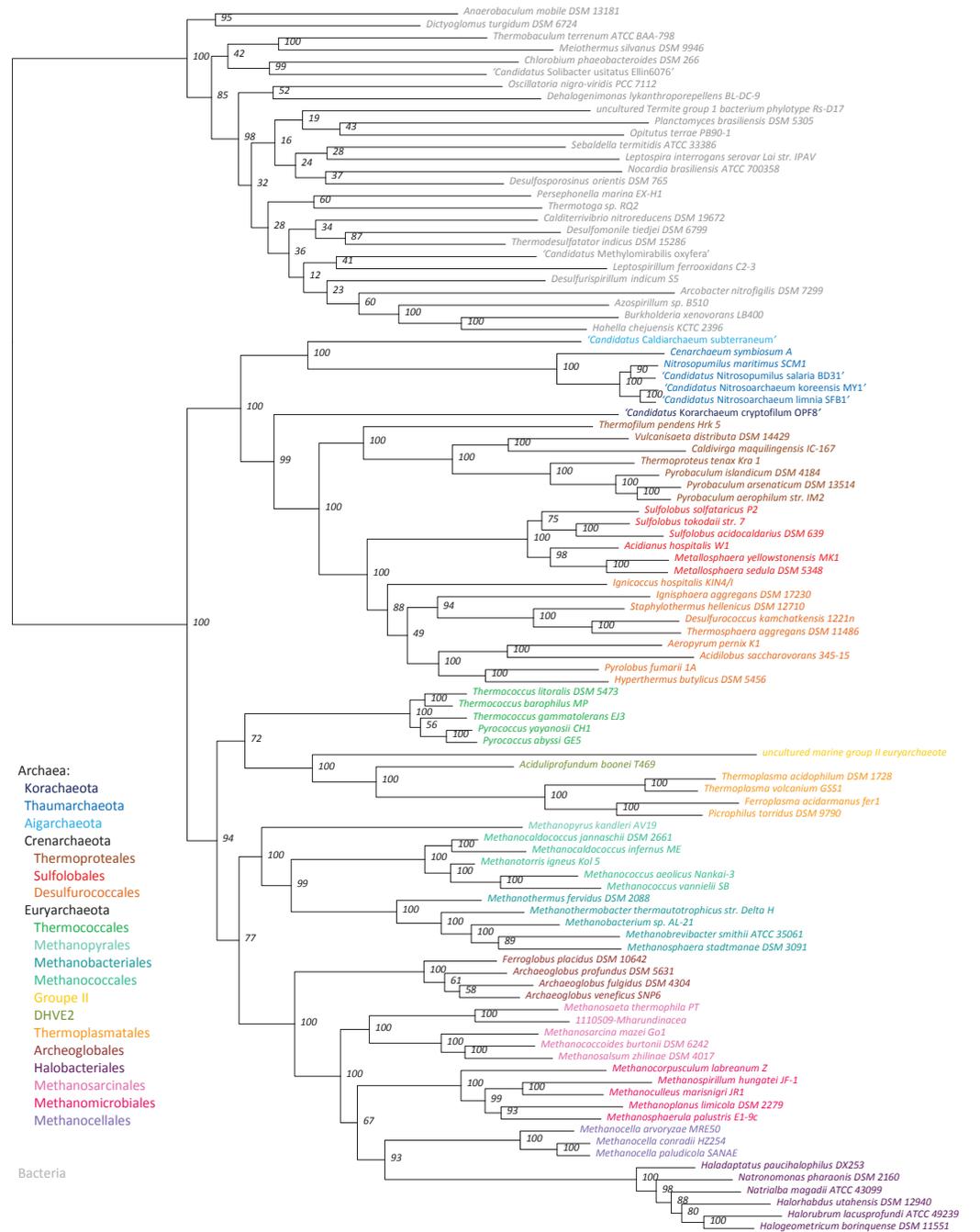
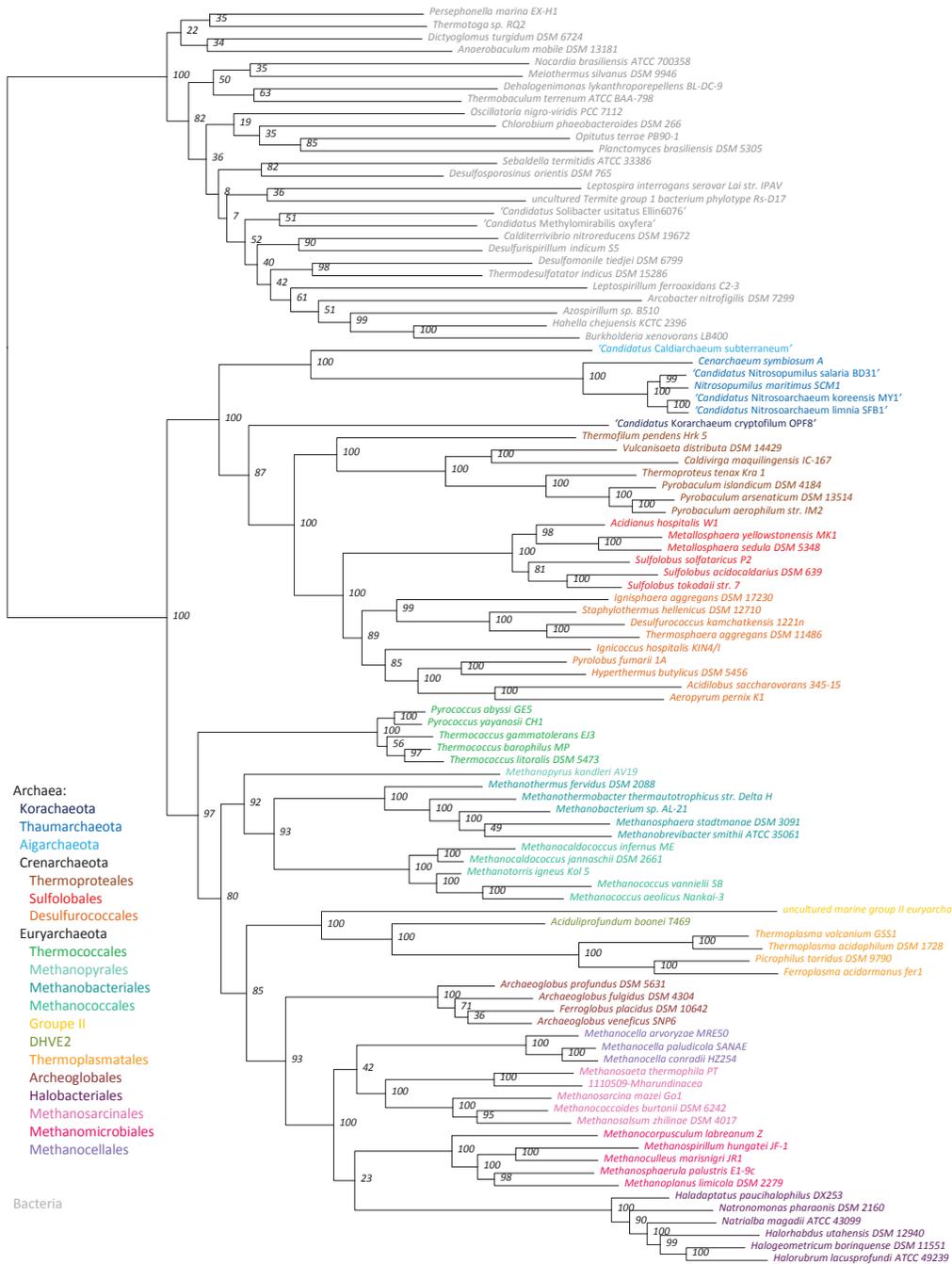


Figure S18
101 species - 5633 positions



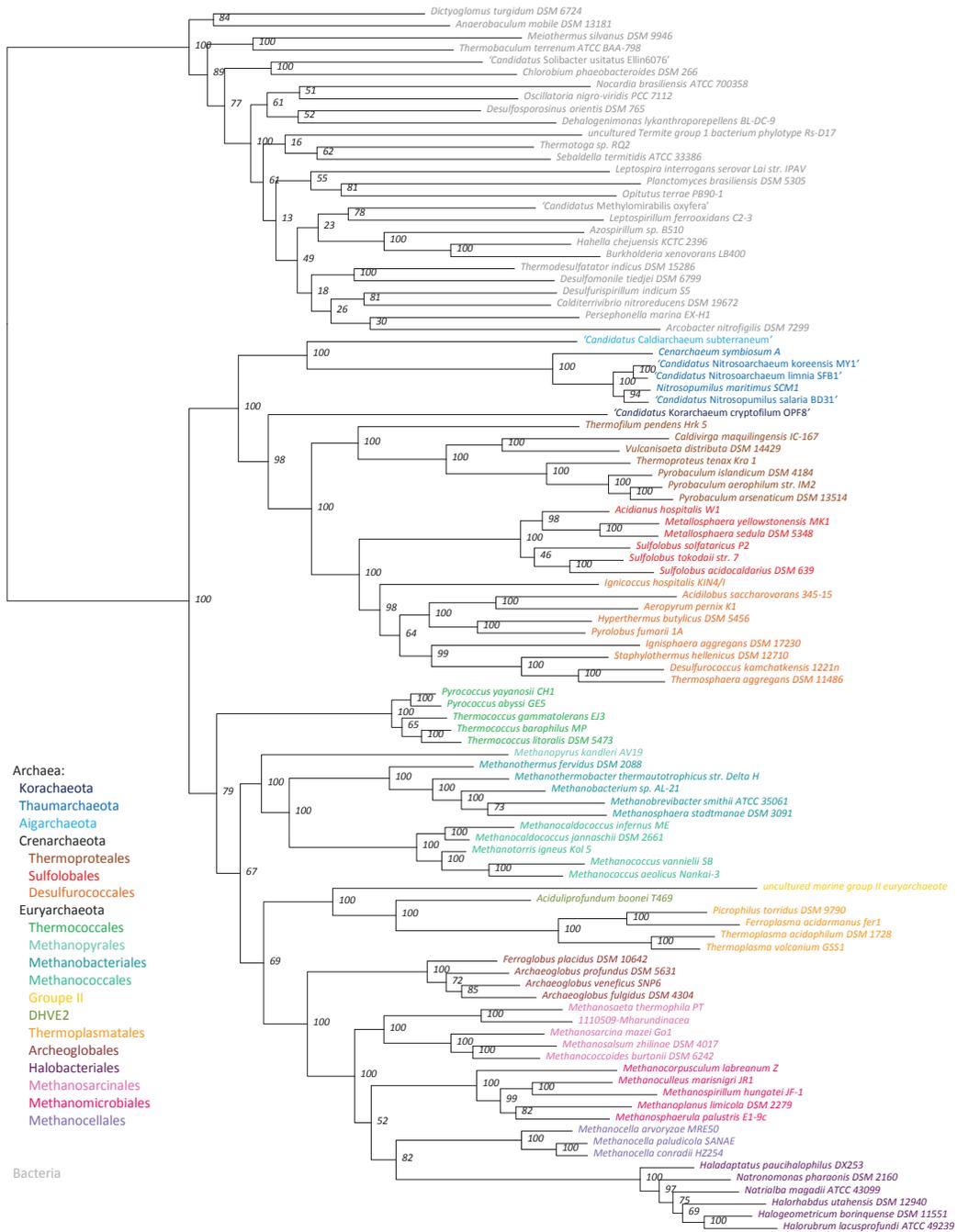


Figure S21
101 species - 8827 positions

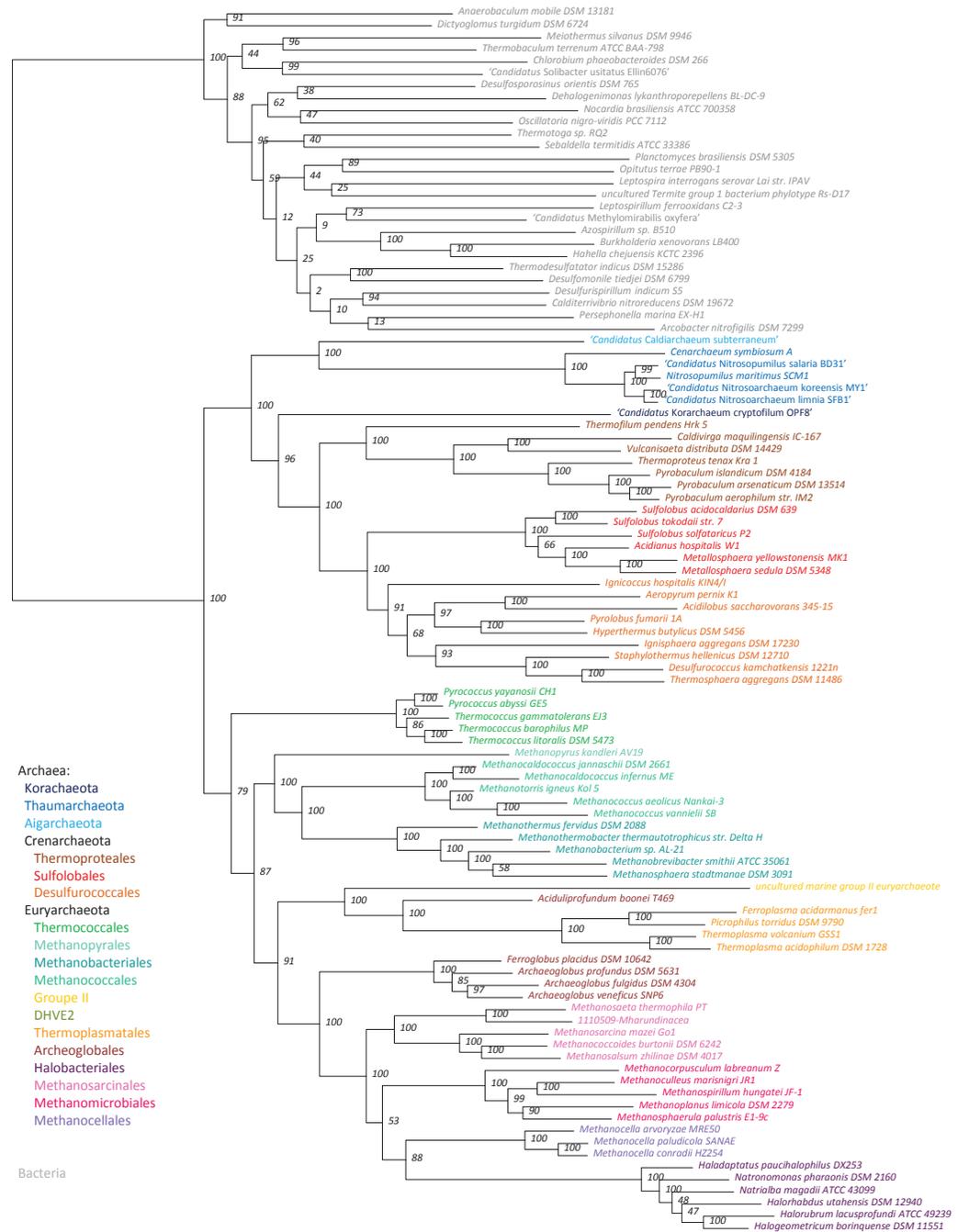


Figure S22
101 species - 9450 positions

Annexe 8 : Liste des 38 nouveaux marqueurs sélectionnés pour l'étude présentée dans le Chapitre 2.

La colonne « Référence thèse » correspond à la référence assignée aux différentes protéines dans cette thèse, utilisée dans les Annexes 9, 10 et 11.

Annotation	Référence thèse	Référence gi	Référence RefSeq
peptidase M50	M004R	161527799	YP_001581625.1
metalloendopeptidase glycoprotease family	M010R	161529041	YP_001582867.1
GTP-binding signal recognition particle	M028R	161528039	YP_001581865.1
cytidyltransferase-like protein	M029R	161528806	YP_001582632.1
beta-lactamase domain-containing protein	M032R	161529047	YP_001582873.1
glutamyl-tRNA(Gln) amidotransferase subunit E	M041R	161527611	YP_001581437.1
tRNA-guanine transglycosylase	M061R	161529044	YP_001582870.1
glutamyl-tRNA(Gln) amidotransferase B subunit	M064R	161528376	YP_001582202.1
exosome complex exonuclease 1	M066R	161527940	YP_001581766.1
MiaB-like tRNA modifying enzyme	M071R	161528769	YP_001582595.1
phenylalanyl-tRNA synthetase alpha subunit	M078R	161528997	YP_001582823.1
argininosuccinate synthase	M138R	161528794	YP_001582620.1
homoserine kinase	M142R	161527553	YP_001581379.1
homoserine dehydrogenase	M143R	161528204	YP_001582030.1
aspartate kinase	M146R	161529262	YP_001583088.1
aspartate carbamoyltransferase	M149R	161529194	YP_001583020.1
phosphoribosylformylglycinamide cyclo-ligase	M151R	161528090	YP_001581916.1
uridylate kinase putative	M153R	161529213	YP_001583039.1
phosphoribosylformylglycinamide synthase II	M156R	161528288	YP_001582114.1
adenylosuccinate lyase	M157R	161528963	YP_001582789.1
RdgB/HAM1 family non-canonical purine NTP pyrophosphatase	M158R	161529038	YP_001582864.1
glutamine amidotransferase class-II	M161R	161528289	YP_001582115.1
phosphoglycerate kinase	M165R	161528008	YP_001581834.1
phosphopantothencysteine decarboxylase/phosphopantothenate--cysteine ligase	M170R	161529229	YP_001583055.1
pyridoxine biosynthesis protein	M171R	161528987	YP_001582813.1
UbiD family decarboxylase	M175R	161529105	YP_001582931.1
porphobilinogen deaminase	M176R	118576540	YP_876283.1
glutamate-1-semialdehyde-2 1-aminomutase	M179R	161527998	YP_001581824.1
GTP1/OBG protein	M195R	161529113	YP_001582939.1
LPPG:FO 2-phospho-L-lactate transferase	M208R	161528134	YP_001581960.1
prolyl-tRNA synthetase	MA17R	315425642	BAJ47301.1
histidyl-tRNA synthetase	MA19R	315425704	BAJ47360.1
signal recognition particle receptor	MA23R	315425933	BAJ47583.1
hydrogenase maturation protein HypF	MA32R	315426449	BAJ48087.1
hydrogenase expression formation protein HypD	MA33R	315426462	BAJ48095.1
5-nucleotidase SurE	MA42R	315426747	BAJ48371.1
conserved hypothetical protein	MA47R	315426923	BAJ48542.1
molybdenum cofactor biosynthesis protein C	MA53R	315427188	BAJ48802.1

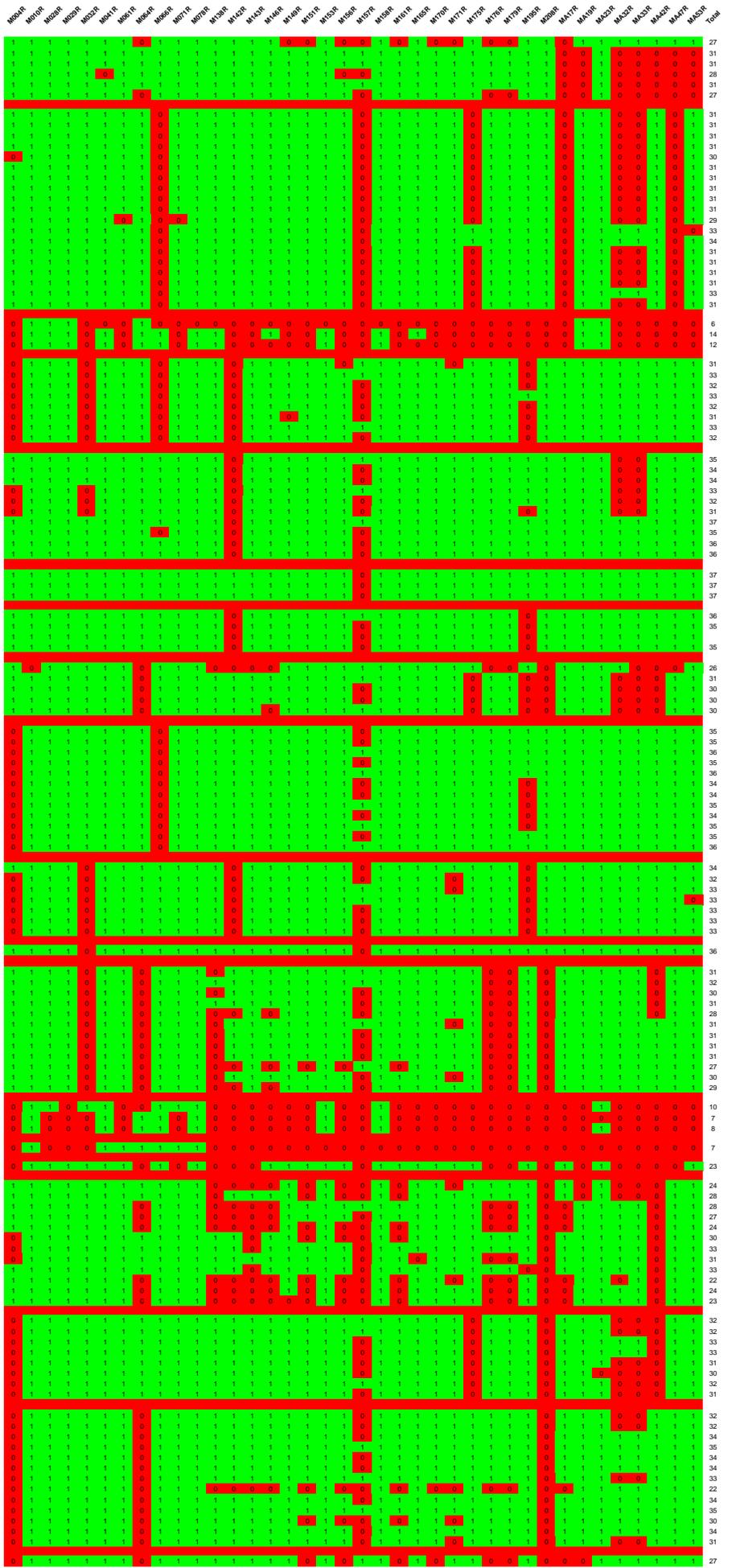
Annexe 9 : Répartition taxonomique des 38 nouveaux marqueurs chez les archées et chez les bactéries.

Chaque ligne correspond à une espèce, chaque colonne à un jeu de données, le croisement entre les deux correspond à la présence (« 1 », case verte), ou à l'absence (« 0 », case rouge) d'homologue de l'espèce dans le jeu de données. La dernière colonne donne le nombre de jeux de données dans lesquels l'espèce est présente.

La liste des archées et des bactéries correspond uniquement aux 129 et 117 génomes utilisés dans ce travail.

Les espèces utilisées pour la construction des jeux de données finaux sont notées par un astérisque.

TAXID	Species	ARCHÉES	Pflanzl	Chlorella	Chlorella	Phylogenetic																							
						MG1	MG2	MG3	MG4	MG5	MG6	MG7	MG8	MG9	MG10	MG11	MG12	MG13	MG14	MG15	MG16	MG17	MG18	MG19	MG20	MG21	MG22	MG23	MG24
1	311458	<i>Candidatus Caldiarchaeum subterraneum</i>	* Algirchaeta	Unclassified	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	27		
2	414004	<i>Cenarchaeum symbiosum A</i>	* Thaumarchaeota	Cenarchaeales	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	31		
3	436308	<i>Nitrosopumilus maritimus DSM1</i>	* Thaumarchaeota	Nitrososummales	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	31		
4	886738	<i>Candidatus Nitrososphaera immitis SFB1</i>	* Thaumarchaeota	Nitrososummales	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	28		
5	1001984	<i>Candidatus Nitrososphaera koreana MT1</i>	* Thaumarchaeota	Nitrososummales	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	31		
6	859350	<i>Candidatus Nitrosopumilus salariae B031</i>	* Thaumarchaeota	Nitrososummales	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	27		



Annexe 10 : Liste des matrices utilisées pour l'analyse en désaturation par sélection de sites pour l'étude exposée dans le Chapitre 2.

Matrice	Nombre de positions
S_5	939
S_{10}	2032
S_{13}	2905
S_{16}	3803
S_{19}	4727
S_{23}	5860
S_{27}	6939
S_{32}	7978
S_{40}	8980
Totale	9450

Annexe 11 : Liste des matrices utilisées pour l'analyse en désaturation par sélection de gènes pour l'étude exposée dans le Chapitre 2.

Pour chaque marqueur, la matrice indiquée est la plus grande dans laquelle ce marqueur a été inclus.

Par exemple, les marqueurs de la matrice DTR21 sont aussi dans les matrices DTR07 et DTR14.

Le nombre de distance calculées correspond au nombre de couples d'espèces pour lesquels une distance a été mesurée, et qui sert de diviseur pour le calcul de la distance moyenne (maximum = 10).

Marqueur	Distance moyenne	Nombre de distances calculées	Matrice	Nombre de positions
s4	0,96050367	22	DTR07	870
s12	0,9814734	22		
s19	1,02444372	22		
M151R	1,10259481	22		
M171R	1,11171562	22		
MA42R	1,12795002	15		
M064R	1,23301883	15		
M158R	1,37887247	22	DTR14	1926
M004R	1,40398118	4		
MA23R	1,43100849	22		
I5	1,4362529	22		
I6	1,44130399	22		
M179R	1,44385348	18		
MA33R	1,44517545	22		
MA32R	1,46667928	22	DTR21	3172
s11	1,47497957	22		
M138R	1,51568416	18		
M176R	1,52872207	18		
I23	1,54398756	22		
M028R	1,55431829	22		
s15	1,56131733	22		
M143R	1,57662369	22	DTR28	4175
M157R	1,58667373	22		
M208R	1,62469153	6		
s8	1,63179175	22		
s3	1,64072369	22		
M195R	1,64326337	22		
M153R	1,66182908	22		
s17	1,67199516	22	DTR35	4902
s7	1,67240291	22		
M078R	1,68421457	22		
I11	1,69115896	22		
I24	1,70010601	22		
I3	1,71632184	22		
M066R	1,73363966	15		
I4	1,73706915	22	DTR42	5633
M170R	1,75706427	22		
I10	1,76798183	22		
I13	1,8150769	22		
I29	1,81841707	22		
I2	1,8385409	22		
M165R	1,86667573	22		

M156R	1,93674863	22	DTR49	7245
s5	1,93978399	22		
MA53R	1,94381957	22		
M142R	1,94681564	15		
M029R	1,94934403	22		
MA17R	1,95739969	15		
M175R	1,95802418	15		
s13	1,96731478	22	DTR56	7881
M071R	1,97633225	22		
MA19R	1,98073354	22		
M032R	1,98385148	12		
M146R	1,98393917	22		
M041R	1,99327988	22		
M149R	2,06404805	22		
s9	2,06570095	22	DTR63	8827
M161R	2,11260774	22		
l1	2,1819755	22		
l10e	2,19431942	22		
M061R	2,2061242	22		
l14	2,22846266	22		
MA47R	2,30049458	22		
l22	2,32519427	22	Totale	9450
M010R	2,35635182	22		
s10	2,4054802	22		
s2	2,56797706	22		
l18	2,5921461	22		
l30	2,88936107	22		
l15	3,06879313	22		

Annexe 12 : Matériels supplémentaires de l'Article 3 (Chapitre 3).

Additional files

Additional file 1: Table showing the taxonomic distribution of DnaJ-Fer, DnaK, DnaJ and GrpE proteins in Archaea. Numbers correspond to accession numbers in the NCBI Genpep database in the 92 complete genomes available in July 2011 and that of, 'Ca. Nitrosoarchaeum koreensis', a thaumarchaeotal genome available more recently. The two divergent sequences of DnaK and GrpE found in *Methanococcus vannielii* SB are underlined.

Additional file 2: Unrooted ML tree of the DnaK protein (136 sequences and 444 positions) inferred with TreeFinder and the LG + Γ 4 model. Numbers at nodes represent bootstrap values and Bayesian posterior probabilities computed with TreeFinder and MrBayes, respectively (only values >50% and 0.5 are shown, dashes indicate the corresponding support is inferior to the threshold, whereas when both supports are inferior to the thresholds no support values are indicated). Archaeal sequences are shown with colours according to their taxonomic classification. The scale bar represents the average number of substitutions per site.

Additional file 3: Unrooted ML tree of the DnaJ protein (102 sequences and 227 positions) inferred with TreeFinder and the LG + Γ 4. Numbers at nodes represent bootstrap values and Bayesian posterior probabilities computed with TreeFinder and MrBayes, respectively (only values >50% and 0.5 are shown, dashes indicate that the corresponding support is inferior to the threshold, whereas when both supports are inferior to the thresholds no support values are indicated). Archaeal sequences are shown with colours according to their taxonomic classification. The scale bar represents the average number of substitutions per site.

Additional file 4: Unrooted ML tree of the GrpE protein (101 sequences and 105 positions) inferred with TreeFinder and the LG + Γ 4. Numbers at nodes represent bootstrap values and Bayesian posterior probabilities computed with TreeFinder and MrBayes, respectively (only values >50% and 0.5 are shown, dashes indicate the corresponding value is inferior to the threshold, whereas when both supports are inferior to the thresholds no support values are indicated). Archaeal sequences are shown with colours according to their taxonomic classification. The scale bar represents the average number of substitutions per site.

TaxID	Organism	Group		DnaJ-Fer	GrpE	DnaK	DnaJ	Cluster
414004	<i>Cenarchaeum symbiosum</i> A	Group I.1a	Thaumarchaeota	YP_875357	YP_876871	YP_876872	YP_876873	EKJ
436308	<i>Nitrosopumilus maritimus</i> SCM1	Group I.1a	Thaumarchaeota	YP_001582358	YP_001581433	YP_001581434	YP_001581435	EKJ
886738	' <i>Candidatus Nitrosoarchaeum limnia</i> SFB1'	Group I.1a	Thaumarchaeota	ZP_08256939	ZP_08256267	ZP_08256266	ZP_08256265	EKJ
1088740	' <i>Candidatus Nitrosoarchaeum korensis</i> '	Group I.1a	Thaumarchaeota	ZP_08668202	ZP_08667092	ZP_08667093	ZP_08667094	EKJ
926571	<i>Nitrososphaera viennensis</i> EN76	Group I.1b	Thaumarchaeota	(unpublished)	NA	NA	NA	NA
497727	' <i>Candidatus Nitrososphaera gargensis</i> '	Group I.1b	Thaumarchaeota	(unpublished)	ADK25989	ADK25988	ADK25990.1	EKJ
311458	' <i>Candidatus Caldichaeum subterraneum</i> '		'Aigarchaeota'		BAJ47057	BAJ47058	BAJ47059	EKJ
64091	<i>Halobacterium</i> sp. NRC-1	Halobacteriales	Euryarchaeota		NP_279548	NP_279546	NP_279545	E-KJ
272569	<i>Haloarcula marismortui</i> ATCC 43049	Halobacteriales	Euryarchaeota		YP_137736	YP_137735	YP_137730	EK---J
348780	<i>Natronomonas pharaonis</i> DSM 2160	Halobacteriales	Euryarchaeota		YP_325771	YP_325772	YP_325835	EK/J
362976	<i>Haloquadratum walsbyi</i> DSM 16790	Halobacteriales	Euryarchaeota		YP_658356	YP_658357	YP_658358	EKJ
416348	<i>Halorubrum lacusprofundi</i> ATCC 49239	Halobacteriales	Euryarchaeota		YP_002565953	YP_002565354	YP_002565353	E / KJ
469382	<i>Halogeometricum borinquense</i> DSM 11551	Halobacteriales	Euryarchaeota		ZP_04000404	ZP_04000406	ZP_04000407	E-KJ
478009	<i>Halobacterium salinarum</i> R1	Halobacteriales	Euryarchaeota		YP_001688639	YP_001688637	YP_001688636	E-KJ
485914	<i>Halomicrobium mukohataei</i> DSM 12286	Halobacteriales	Euryarchaeota		YP_003178657	YP_003178658	YP_003178662	EK---J
519442	<i>Halorhabdus utahensis</i> DSM 12940	Halobacteriales	Euryarchaeota		YP_003129229	YP_003129230	YP_003129211	EK / J
543526	<i>Haloterrigena turkmenica</i> DSM 5511	Halobacteriales	Euryarchaeota		YP_003403571	YP_003403570	YP_003403568	E-KJ
547559	<i>Natrialba magadii</i> ATCC 43099	Halobacteriales	Euryarchaeota		YP_003478805	YP_003478720	YP_003478719	E / KJ
889948	' <i>Candidatus Nanosalina</i> sp. J07AB43'	Nanoarchaeota	Euryarchaeota		EGQ43817	EGQ43816	EGQ43815	EKJ
889962	' <i>Candidatus Nanosalinarum</i> sp. J07AB56'	Nanoarchaeota	Euryarchaeota		EGQ40899	EGQ40898	EGQ40051	EK / J
323259	<i>Methanospirillum hungatei</i> JF-1	Methanomicrobiales	Euryarchaeota		YP_501625	YP_501624	YP_501623	EKJ
368407	<i>Methanoculleus marisnigri</i> JR1	Methanomicrobiales	Euryarchaeota		YP_001046977	YP_001046976	YP_001046975	EKJ
410358	<i>Methanococcus profundus</i> ATCC 35061	Methanomicrobiales	Euryarchaeota		YP_001030381	YP_001030380	YP_001030379	EKJ
456442	' <i>Candidatus Methanoregula boonei</i> 6A8'	Methanomicrobiales	Euryarchaeota		YP_001404367	YP_001404368	YP_001404369	EKJ
456442	' <i>Candidatus Methanoregula boonei</i> 6A8'	Methanomicrobiales	Euryarchaeota		YP_001404765	YP_001404766	YP_001404767	EKJ
521011	<i>Methanosphaerula palustris</i> E1-9c	Methanomicrobiales	Euryarchaeota		YP_002466649	YP_002466650	YP_002466651	EKJ
521011	<i>Methanosphaerula palustris</i> E1-9c	Methanomicrobiales	Euryarchaeota		YP_002467356	YP_002467355	YP_002467354	EKJ
304371	<i>Methanocella paludicola</i> SANA	Methanocellales	Euryarchaeota		YP_003357075	YP_003357074	YP_003357073	EKJ
351160	uncultured methanogenic archaeon RC-1	Methanocellales	Euryarchaeota		YP_684746	YP_684745	YP_684744	EKJ
188937	<i>Methanosarcina acetivorans</i> C2A	Methanosarcinales	Euryarchaeota		NP_616411	NP_616412	NP_616413	EKJ
192952	<i>Methanosarcina mazelii</i> Go1 (NP_634529)	Methanosarcinales	Euryarchaeota		NP_634530	NP_634529	NP_634528	EKJ
259564	<i>Methanococcoides burtonii</i> DSM 6242	Methanosarcinales	Euryarchaeota		YP_565987	YP_565980	YP_565979	EKJ
269797	<i>Methanosarcina barkeri</i> str. <i>Fusaro</i>	Methanosarcinales	Euryarchaeota		YP_306881	YP_306886	YP_306885	EKJ
349307	<i>Methanosarcina thermophila</i> PT	Methanosarcinales	Euryarchaeota		YP_843162	YP_843161	YP_843160	EKJ
547558	<i>Methanohalophilus mahii</i> DSM 5219	Methanosarcinales	Euryarchaeota		YP_003542758	YP_003542757	YP_003542756	EKJ
224325	<i>Archaeoglobus fulgidus</i> DSM 4304	Archaeoglobales	Euryarchaeota					
572546	<i>Archaeoglobus profundus</i> DSM 5631	Archaeoglobales	Euryarchaeota					
589924	<i>Ferroglobus placidus</i> DSM 10642	Archaeoglobales	Euryarchaeota					
439481	<i>Aciduliprofundum boonei</i> T469 (YP_003483718)	DHEV2	Euryarchaeota		YP_003483717.1	YP_003483718	YP_003483720.1	EK-J
263820	<i>Picrophilus torridus</i> DSM 9790	Thermoplasmatales	Euryarchaeota		YP_023617	YP_023618	YP_023619	EKJ
273075	<i>Thermoplasma acidophilum</i> DSM 1728	Thermoplasmatales	Euryarchaeota		NP_394545	NP_394546	NP_394547	EKJ
273116	<i>Thermoplasma volcanium</i> GSS1	Thermoplasmatales	Euryarchaeota		NP_111008	NP_111007	NP_111006	EKJ
333146	<i>Ferroplasma acidimanus</i> fer1	Thermoplasmatales	Euryarchaeota		ZP_05571351	ZP_05571352	ZP_05571353	EKJ
425595	' <i>Candidatus Micrarchaeum acidiphilum</i> ARMAN-2'	ARMAN	Euryarchaeota		EET90013	EET90012	EET90011	EKJ
662760	' <i>Candidatus Parvarchaeum acidiphilum</i> ARMAN-4'	ARMAN	Euryarchaeota		EEZ92631	EEZ93272	EEZ93273	E / KJ
662762	' <i>Candidatus Parvarchaeum acidiphilum</i> ARMAN-5'	ARMAN	Euryarchaeota		EPF92431	EPF92422	EPF92429	EKJ
243232	<i>Methanocaldococcus jannaschii</i> DSM 2661	Methanococcales	Euryarchaeota					
267377	<i>Methanococcus maripaludis</i> S2	Methanococcales	Euryarchaeota					
402880	<i>Methanococcus maripaludis</i> C5	Methanococcales	Euryarchaeota					
406327	<i>Methanococcus vanielii</i> SB	Methanococcales	Euryarchaeota		YP_001322618	YP_001322619		EK
419665	<i>Methanococcus aeolicus</i> Nankai-3	Methanococcales	Euryarchaeota					
426368	<i>Methanococcus maripaludis</i> C7	Methanococcales	Euryarchaeota					
444158	<i>Methanococcus maripaludis</i> C6	Methanococcales	Euryarchaeota					
456320	<i>Methanococcus voltae</i> A3	Methanococcales	Euryarchaeota					
573063	<i>Methanocaldococcus infernus</i> ME	Methanococcales	Euryarchaeota					
573064	<i>Methanocaldococcus fervens</i> AG86	Methanococcales	Euryarchaeota					
579137	<i>Methanocaldococcus vulcanius</i> M7	Methanococcales	Euryarchaeota					
187420	<i>Methanothermobacter thermautotrophicus</i> str. <i>Delta</i> H	Methanobacteriales	Euryarchaeota		NP_276410	NP_276411	NP_276412	EKJ
339860	<i>Methanospaera stadmanae</i> DSM 3091	Methanobacteriales	Euryarchaeota		YP_448530	YP_448529	YP_448528	EKJ
420247	<i>Methanobrevibacter smithii</i> ATCC 35061	Methanobacteriales	Euryarchaeota		YP_001273681	YP_001273682	YP_001273683	EKJ
634498	<i>Methanobrevibacter ruminantium</i> M1	Methanobacteriales	Euryarchaeota		YP_003424780	YP_003424781	YP_003424782	EKJ
190192	<i>Methanopyrus kandleri</i> AV19	Methanopyrales	Euryarchaeota					
69014	<i>Thermococcus kodakarensis</i> KOD1	Thermococcales	Euryarchaeota					
70601	<i>Pyrococcus horikoshii</i> OT3	Thermococcales	Euryarchaeota					
186497	<i>Pyrococcus furiosus</i> DSM 3638	Thermococcales	Euryarchaeota					
246969	<i>Thermococcus</i> sp. AM4	Thermococcales	Euryarchaeota					
272844	<i>Pyrococcus abyssi</i> GE5	Thermococcales	Euryarchaeota					
523850	<i>Thermococcus onnurineus</i> NA1	Thermococcales	Euryarchaeota					
593117	<i>Thermococcus gammatolerans</i> EJ3	Thermococcales	Euryarchaeota					
604354	<i>Thermococcus sibiricus</i> MM 739	Thermococcales	Euryarchaeota					
228908	<i>Nanoarchaeum equitans</i> Kin4-M	Nanoarchaeota	Euryarchaeota					
272557	<i>Aeropyrum pernix</i> K1	Desulfurococcales	Crenarchaeota					
399550	<i>Staphylothermus marinus</i> F1	Desulfurococcales	Crenarchaeota					
415426	<i>Hyperthermus butylicus</i> DSM 5456	Desulfurococcales	Crenarchaeota					
453591	<i>Ignicoccus hospitalis</i> KIN4/I	Desulfurococcales	Crenarchaeota					
490899	<i>Desulfurococcus kamchatkensis</i> 1221n	Desulfurococcales	Crenarchaeota					
591019	<i>Staphylothermus hellenicus</i> DSM 12710	Desulfurococcales	Crenarchaeota					
633148	<i>Thermosphaera aggregans</i> DSM 11486	Desulfurococcales	Crenarchaeota					
273057	<i>Sulfolobus solfatarius</i> P2	Sulfolobales	Crenarchaeota					
273063	<i>Sulfolobus tokodaii</i> str. 7	Sulfolobales	Crenarchaeota					
330779	<i>Sulfolobus acidocaldarius</i> DSM 639	Sulfolobales	Crenarchaeota					
399549	<i>Metallosphaera sedula</i> DSM 5348	Sulfolobales	Crenarchaeota					
419942	<i>Sulfolobus islandicus</i> Y.N.15.51	Sulfolobales	Crenarchaeota					
425944	<i>Sulfolobus islandicus</i> L.D.8.5	Sulfolobales	Crenarchaeota					
426118	<i>Sulfolobus islandicus</i> M.16.4	Sulfolobales	Crenarchaeota					
427317	<i>Sulfolobus islandicus</i> M.14.25	Sulfolobales	Crenarchaeota					
427318	<i>Sulfolobus islandicus</i> M.16.27	Sulfolobales	Crenarchaeota					
429572	<i>Sulfolobus islandicus</i> L.S.2.15	Sulfolobales	Crenarchaeota					
439386	<i>Sulfolobus islandicus</i> Y.G.57.14	Sulfolobales	Crenarchaeota					
178306	<i>Pyrobaculum aerophilum</i> str. IM2	Thermoproteales	Crenarchaeota					
340102	<i>Pyrobaculum arsenaticum</i> DSM 13514	Thermoproteales	Crenarchaeota					
368408	<i>Thermofilum pendens</i> Hrk 5	Thermoproteales	Crenarchaeota					
384616	<i>Pyrobaculum islandicum</i> DSM 4184	Thermoproteales	Crenarchaeota					
397948	<i>Caldivirus maquilingensis</i> IC-167	Thermoproteales	Crenarchaeota					
410359	<i>Pyrobaculum calidifontis</i> JCM 11548	Thermoproteales	Crenarchaeota					
444157	<i>Thermoproteus neutrophilus</i> V24St4	Thermoproteales	Crenarchaeota					
374847	' <i>Candidatus Korarchaeum cryptoflum</i> OPF8'		Korarchaeota					



(Petitjean et al., Suppl. Figure S1)

