

Title : Streamlining and large ancestral genomes in Archaea inferred with a phylogenetic birth-and-death model

Authors : Miklós Csűrös¹ and István Miklós²

Authors' affiliations: ¹ Department of Computer Science and Operations Research, University of Montréal, Canada. ² Rényi Institute of Mathematics, Hungarian Academy of Sciences, Budapest, Hungary.

Corresponding author: Miklós Csűrös. Département d'informatique et de recherche opérationnelle, Université de Montréal, C.P. 6128, succursale Centre-Ville, Montréal, QC, H3C 3J7, Canada. Tel: +1 514 343-6111 extension 1655. Fax: +1 514 343-6111. E-mail: csuros@iro.umontreal.ca.

Abstract

1
2 Homologous genes originate from a common ancestor through vertical
3 inheritance, duplication or horizontal gene transfer. Entire homolog fam-
4 ilies spawned by a single ancestral gene can be identified across multiple
5 genomes based on protein sequence similarity. The sequences, however, do
6 not always reveal conclusively the history of large families. In order to study
7 the evolution of complete gene repertoires, we propose here a mathematical
8 framework that does not rely on resolved histories. We show that so-called
9 phylogenetic profiles, formed by family sizes across multiple genomes, are
10 sufficient to infer principal evolutionary trends. The main novelty in our ap-
11 proach is an efficient algorithm to compute the likelihood of a phylogenetic
12 profile in a model of birth-and-death processes acting on a phylogeny.

13 We examine known gene families in 28 archaeal genomes using a proba-
14 bilistic model that involves lineage- and family-specific components of gene
15 acquisition, duplication, and loss. The model enables us to consider all pos-
16 sible histories when inferring statistics about archaeal evolution. According
17 to our reconstruction, most lineages are characterized by a net *loss* of gene
18 families. Major increases in gene repertoire have occurred only a few times.
19 Our reconstruction underlines the importance of persistent streamlining pro-
20 cesses in shaping genome composition in Archaea. It also suggests that early
21 archaeal genomes were as complex as typical modern ones, and even show
22 signs, in the case of the methanogenic ancestor, of an extremely large gene
23 repertoire.

24 **Introduction**

25 The evolution of homologous gene families, i.e., genes of common ancestry, is
26 enmeshed within species histories in a complex manner (Koonin, 2005). Con-
27 comitantly with the diversification of organismal lineages, gene families expand
28 by duplications, individual genes get eliminated, and new genes arrive by lateral
29 transfer. It is now clear that *de novo* gene formation and vertical processes (Snel
30 *et al.*, 2002; Henikoff *et al.*, 1997), such as duplication and loss, act in concert with
31 horizontal gene transfer (Boucher *et al.*, 2003; Gogarten and Townsend, 2005).

32 Gene families are identified in current practice by pairwise sequence compar-
33 isons, coupled with the clustering of postulated homolog pairs (Tatusov *et al.*,
34 1997; Alexeyenko *et al.*, 2006) The phylogenetic profile of a gene family com-
35 prises the family size across a set of organisms, i.e., the number of homologs within
36 the same family in each genome. Such profiles are extremely informative even
37 without taking the gene sequences into account: profile data sets have been used
38 to construct organismal phylogenies (Fitz-Gibbon and House, 1999; Snel *et al.*,
39 1999; Tekaiia *et al.*, 1999) and to infer ancestral gene content (Mirkin *et al.*, 2003;
40 Iwasaki and Takagi, 2007); similar and complementary profiles hint at functional
41 associations (Tatusov *et al.*, 1997; Pellegrini *et al.*, 1999). Considering various
42 evolutionary processes in a mathematical model of gene family evolution is chal-
43 lenging. One main element that distinguishes the present study from past work
44 is the elaboration of a likelihood framework for phylogenetic profiles that simul-
45 taneously accounts for gene duplication, loss, and acquisition. In particular, we
46 describe an algorithm for the exact computation of the likelihood in a phylogenetic
47 gain-loss-duplication model.

48 The present study uses a gain-loss-duplication model to address gene content
49 evolution in Archaea. Relying on a complete set of known homolog families in 28
50 sequenced genomes, we inferred lineage- and family-specific statistics. In a pre-
51 cursory step, we constructed a plausible phylogeny using 88 universally conserved
52 proteins, which we believe is a noteworthy result on its own, as the phylogeny
53 resolves some problematic euryarchaeal branching orders (involving Thermoplas-
54 matales, Methanopyrus and Methanobacteriales) confidently. Gene loss emerges
55 in our analysis as the dominant force that has shaped archaeal genomes through-
56 out their history. Apparently, genome streamlining has been an ongoing process
57 in all lineages with a fairly constant intensity, apart from dramatic genome com-
58 pactions in endosymbiotic Archaea. Our reconstruction suggests that early Ar-
59 chaea had a comparable genomic complexity to today's organisms. In particular,
60 the euryarchaeal ancestor of two classes of methanogens had a very large genome,
61 resulting from one of the rare upsurges in gene content, similarly to some modern
62 lineages of Methanosarcina and Halobacteria.

63 **Methods**

64 **Phylogenetic profiles in Archaea**

65 Phylogenetic profiles, sequences, and functional annotations were downloaded from
66 the arCOG database of orthologous gene clusters in Archaea (Makarova *et al.*,
67 2007) at <ftp://ftp.ncbi.nih.gov/pub/wolf/COGs/arCOG>. The pro-
68 files were amended with data on lineage-specific singletons and inparalog families
69 that have no archaeal homologs outside of one genome (Yuri Wolf, personal com-
70 munication), which was produced in the process of compiling the arCOG database.

71 The following organisms are included in the study: *Archaeoglobus fulgidus*
72 (Arcfu), *Haloarcula marismortui* ATCC 43049 (Halma), *Halobacterium* sp. strain
73 NRC-1 (Halsp), *Methanosarcina acetivorans* (Metac), *Methanococcoides burtonii*
74 DSM 6242 (Metbu), *Methanoculleus marisnigri* JR1 (Metcu), *Methanospirillum*
75 *hungatei* JF-1 (Methu), *Methanocaldococcus jannaschii* (Metja), *Methanopyrus*
76 *kandleri* (Metka), *Methanosarcina mazei* (Metma), *Methanococcus maripaludis*
77 S2 (Metmp), *Methanosphaera stadtmanae* (Metst), *Methanothermobacter thermoau-*
78 *totrophicus* (Metth), *Nanoarchaeum equitans* (Naneq), *Picrophilus torridus* DSM
79 9790 (Picto), *Pyrococcus abyssi* (Pyrab), *Pyrococcus furiosus* (Pyrfu), *Thermo-*
80 *plasma acidophilum* (Theac), *Thermococcus kodakaraensis* KOD1 (Theko), *Ther-*
81 *moplasma volcanium* (Thevo), *Aeropyrum pernix* (Aerpe), *Caldivirga maquilin-*
82 *gensis* IC-167 (Calma), *Cenarchaeum symbiosum* (Censy), *Hyperthermus butylicus*
83 (Hypbu), *Pyrobaculum aërophilum* (Pyræ), *Sulfolobus solfataricus* (Sulso); *Sul-*
84 *folobus acidocaldarius* DSM 639 (Sulac), *Thermofilum pendens* Hrk 5 (Thepe)
85 with the last eight classified as crenarchaeota. The abbreviations are those used
86 by (Makarova *et al.*, 2007) and the arCOG database.

87 **Reconstruction of archaeal phylogeny**

88 The phylogeny was constructed using concatenated multiple alignments of se-
89 lected orthologous protein sequences. The sequences were chosen from the arCOG
90 database based on phylogenetic profiles: we selected all arCOG groups where ev-
91 ery studied genome contained exactly one homolog. There are 88 such groups
92 (see Supplemental Material for sequences), and 46 of those correspond to ribo-
93 somal proteins. Alignments were done using the program Muscle (Edgar, 2004).
94 Phylogenies were built by likelihood maximization using PhyML (Guindon and
95 Gascuel, 2003), with the Jones-Taylor-Thornton substitution model and eight dis-
96 crete Gamma categories and invariant sites. The expected number of substitutions
97 per amino acid site was computed on each edge for the ribosomal proteins in the
98 JTT+I+ Γ 8 model by PhyML. Bootstrap support values for the branches were com-
99 puted by PhyML, using 500 replicates.

100 **Inference of gene content evolution**

101 We maximized the likelihood (see below for the likelihood computation) of the
102 data set using a gain-loss-duplication model with a Poisson distribution at the root
103 and four discrete Gamma categories capturing rate variation across families, for
104 edge length t_f and duplication λ_f each. For a given set of model parameters (three
105 parameters — $\hat{t}_e \hat{\kappa}_e$, $\hat{t}_e \hat{\mu}_e$, $\hat{t}_e \hat{\lambda}_e$ — per edge, one for the root's Poisson param-
106 eter Γ , and two Gamma shape parameters for rate variation), the likelihood of
107 each family was computed using (1) with the described methods of manipulating
108 rate variation and correcting for absent profiles. The data set's likelihood (i.e.,
109 the product of family likelihoods) was then maximized numerically as a function

110 of the model parameters, using custom-made software implementing the Broyden-
 111 Fletcher-Goldfarb-Shanno conjugate gradient method and Brent’s one-dimensional
 112 optimization method (Press *et al.*, 1997). Family sizes and lineage-specific events
 113 (gains, losses, expansions, contractions) were computed using posterior probabilities
 114 in the optimized gain-loss-duplication model.

115 **Phylogenetic birth-and-death model**

116 A *phylogenetic birth-and-death model* formalizes the evolution of an organism-
 117 specific census variable along a rooted phylogeny T . We consider only binary
 118 phylogenies here; the full set of methods applicable to multi-furcating phylogenies
 119 is described in the Supporting Information. The model specifies edge lengths, as
 120 well as birth-and-death processes (Ross, 1996; Kendall, 1949) acting on the edges.
 121 Populations of identical individuals evolve along the tree from the root towards
 122 the leaves by Galton-Watson processes. At non-leaf nodes of the tree, populations
 123 are instantaneously copied to evolve independently along the adjoining descen-
 124 dant edges. Let the random variable $\xi(x) \in \{0, 1, 2, \dots\}$ denote the population
 125 count at every node $x \in \mathcal{V}(T)$. Every edge xy is characterized by a loss rate μ_{xy} ,
 126 a duplication rate λ_{xy} and a gain rate κ_{xy} . If $(X(t): t \geq 0)$ is a linear birth-
 127 and-death process (Kendall, 1949; Takács, 1962) with these rate parameters, then
 128 $\mathbb{P}\{\xi(y) = m \mid \xi(x) = n\} = \mathbb{P}\{X(t_{xy}) = m \mid X(0) = n\}$, where $t_{xy} > 0$ is
 129 the edge length, which defines the time interval during which the birth-and-death
 130 process runs. The joint distribution of $(\xi(x): x \in \mathcal{V}(T))$ is determined by the phy-
 131 logeny, the edge lengths and rates, along with the distribution at the root ρ , denoted
 132 as $\gamma(n) = \mathbb{P}\{\xi(\rho) = n\}$.

133 It is assumed that one can observe the population counts at the terminal nodes

134 (i.e., leaves), but not at the inner nodes of the phylogeny. Since individuals are
 135 considered identical, we are also ignorant of the ancestral relationships between in-
 136 dividuals within and across populations. The population counts at the leaves form
 137 a *phylogenetic profile*, which is formally a function $\Phi: \mathcal{L}(T) \mapsto \{0, 1, 2, \dots\}$,
 138 where $\mathcal{L}(T) \subset \mathcal{V}(T)$ denote the set of leaf nodes. Our central problem is to com-
 139 pute the likelihood of a profile, i.e., the probability of the observed counts for fixed
 140 model parameters. Define the notation $\Phi(\mathcal{L}') = (\Phi(x): x \in \mathcal{L}')$ for the partial
 141 profile within a subset $\mathcal{L}' \subseteq \mathcal{L}(T)$. Similarly, let $\xi(\mathcal{L}') = (\xi(x): x \in \mathcal{L}')$ denote
 142 the vector-valued random variable composed of individual population counts. The
 143 *likelihood* of Φ is the probability $L = \mathbb{P}\{\xi(\mathcal{L}(T)) = \Phi\}$. Let T_x denote the sub-
 144 tree of T rooted at node x . Define the *survival count range* M_x for every node x as
 145 $M_x = \sum_{y \in \mathcal{L}(T_x)} \Phi(y)$. Clearly, the ranges can be calculated easily in a postorder
 146 traversal.

147 For our discussion, we borrow standard terminology applied to homologous
 148 genes (Sonnhammer and Koonin, 2002). For every edge xy , the population of
 149 node y can be split by ancestry at node x : *inparalog* groups are formed by the
 150 progenies of each individual at x and a *xenolog* group is formed by the individuals
 151 whose ancestor immigrated into the population. When $\xi(x) = n$ on the edge xy ,
 152 then $\xi(y) = \eta + \sum_{i=1}^n \zeta_i$, where η is the xenolog group size, and ζ_i are the indepen-
 153 dent and identically distributed inparalog group sizes. The distribution of xenolog
 154 and inparalog group sizes is the well-characterized transient distribution of the ap-
 155 propriate linear birth-and-death processes (Karlin and McGregor, 1958; Kendall,
 156 1949; Takács, 1962); see Supplemental Material. Namely, each ζ_i has a shifted
 157 geometric distribution, and for $\kappa > 0$, η has a negative binomial or Poisson distri-
 158 bution. The distributions' parameters are known functions of the edge length t_{xy}

159 and rates $\kappa_{xy}, \lambda_{xy}, \mu_{xy}$.

160 **Surviving lineages**

161 A key factor in inferring the likelihood formulas is the probability that a given in-
162 dividual at a tree node x has no descendants at the leaves within the subtree rooted
163 at x . The corresponding *extinction probability* is denoted by D_x , which can be
164 computed in a postorder traversal (Csűrös and Miklós, 2006). An individual at
165 node x is referred to as *surviving* if it has at least one progeny at the leaves de-
166 scending from x . Let $\Xi(x)$ denote the number of surviving individuals at each
167 node x . The number of surviving xenologs and inparalogs follow the same class of
168 distributions as the total number of xenologs and inparalogs (see Supplemental Ma-
169 terial). Consequently, if $\xi(x) = n$ on edge xy , then $\Xi(y) = \eta + \sum_{i=1}^n \zeta_i$, where η
170 is the surviving xenolog count with a Poisson or negative binomial distribution, and
171 ζ_i are surviving paralog counts, with negative binomial distributions. The distri-
172 butions' parameters can be computed explicitly using the process parameters and
173 the extinction probabilities. In the formulas to follow, we use the probabilities
174 $w_y^*[m|n] = \mathbb{P}\{\eta + \sum_{i=1}^n \zeta_i = m; \forall \zeta_i > 0\}$, which can be computed by dynamic
175 programming for all $n, m \leq M_y$ in $O(M_y^2)$ time (see Supplemental Material).

176 **Computing the likelihood**

177 We compute the likelihood using *conditional survival likelihoods* defined as the
178 probability of observing the partial profile within T_x given the number of surviving
179 individuals $\Xi(x)$: $L_x[n] = \mathbb{P}\left\{\xi(\mathcal{L}(T_x)) = \Phi(\mathcal{L}(T_x)) \mid \Xi(x) = n\right\}$. For $m >$
180 M_x , $L_x[m] = 0$. For values $m = 0, 1, \dots, M_x$, the conditional survival likelihoods
181 can be computed recursively as shown below.

182 If node x is a leaf, then

$$L_x[n] = \begin{cases} 0 & \text{if } n \neq \Phi(x); \\ 1 & \text{if } n = \Phi(x). \end{cases}$$

If x is an inner node with children x_1, x_2 , then $L_x[n]$ can be expressed using $L_{x_i}[\cdot]$ and auxiliary values $B_{i;\cdot}$, for $i = 1, 2$ in the following manner. Auxiliary values $B_{i;t,s}$ are defined for $i = 1, 2$ and $s = 0, \dots, M_{x_i}$ as follows.

$$B_{i;0,s} = \sum_{m=0}^{M_{x_i}} w_{x_i}^*[m|s] L_{x_i}[m] \quad \{0 \leq s \leq M_{x_i}\}$$

$$B_{2;t,M_{x_2}} = G_{x_2}(0) B_{2;t-1,M_{x_2}}$$

$$B_{2;t,s} = B_{2;t-1,s+1} + G_{x_2}(0) B_{2;t-1,s} \quad \{0 \leq s < M_{x_2}\}$$

183 where $G_{x_i}(k) = \mathbb{P}\{\zeta = k\}$ for a surviving inparalog group at x_i . In the above
184 equations, $0 < t \leq M_{x_1}$. For all $n = 0, \dots, M_x$

$$L_x[n] = (1 - D_x)^{-n} \sum_{\substack{0 \leq t \leq M_{x_1} \\ 0 \leq s \leq M_{x_2} \\ t+s=n}} \binom{n}{s} (D_{x_1})^s B_{1;0,t} B_{2;t,s}.$$

185 The complete likelihood is computed as

$$L = \sum_{m=0}^{M_\rho} L_\rho[m] \mathbb{P}\{\Xi(\rho) = m\}.$$

186 For some parametric distributions γ , there is a closed formula for $\mathbb{P}\{\Xi(\rho) = m\}$. In
187 particular, if γ is the stationary distribution for a gain-loss-duplication or a gain-loss

188 models, then $\Xi(\rho)$ has a negative binomial or Poisson distribution, respectively.

189 The likelihood for a Poisson distribution at the root is

$$L = \sum_{m=0}^{M_\rho} L_\rho[m] \exp\left(-\Gamma(1 - D_\rho)\right) \frac{(\Gamma(1 - D_\rho))^m}{m!} \quad (1)$$

190 where Γ is the mean family size at the root.

191 The likelihood formula (1) is corrected in order to account for the fact that the
192 data set does not contain all-absent profiles with $\Phi(x) = 0$ for all leaves x , in a
193 manner analogous to (Felsenstein, 1992).

194 Family-specific rate variation is considered by computing the likelihood val-
195 ues for each discrete rate category c characterized by factors $(t_c, \kappa_c, \mu_c, \lambda_c)$. The
196 factors in our analysis are either constant 1, or correspond to the expected values
197 within the four quartiles of a Gamma distribution with mean 1.

198 **Results and discussion**

199 **Computational analysis of phylogenetic profiles**

200 Birth-and-death processes are commonly used to model a population of identi-
201 cal individuals (Kendall, 1949; Karlin and McGregor, 1958) and waiting queues
202 (Takács, 1962). Their use in modeling gene family evolution is justified by the
203 fact that losses and duplications seem to occur independently between the mem-
204 bers of multi-gene families (Nei and Rooney, 2005). The most general process we
205 consider is a gain-loss-duplication process which is characterized by the rates of
206 gain κ , loss μ and duplication λ : a population of size n grows by a rate of $(\lambda n + \kappa)$
207 and decreases by a rate of μn . In our context, the population comprises homologs
208 of a given family in the genome. Gene acquisition occurs with a rate of κ , combin-
209 ing various means such as innovation and lateral transfer. We model gene family
210 evolution in a phylogenetic setting by associating gain-loss-duplication processes
211 with the branches of a phylogenetic tree. The corresponding phylogenetic birth-
212 and-death model defines a probabilistic framework for the evolution of gene family
213 size. The observed family sizes at the terminal nodes form a phylogenetic profile.
214 In principle, a phylogenetic birth-and-death model suits likelihood-based inference
215 since it is a probabilistic graphical model (Jordan, 2004) with a tree structure. The
216 mathematical difficulties stem from the fact that the state space of the processes
217 (i.e., family size) is infinitely large. Consequently, routine computational tech-
218 niques used to analyze molecular sequence evolution (Felsenstein, 1981) are not
219 applicable. Previously proposed likelihood methods (Hahn *et al.*, 2005; Spencer
220 *et al.*, 2006; Iwasaki and Takagi, 2007) have sidestepped the infinity problem by

221 using approximative calculations with bounds on maximal family size.

222 We have introduced (Csűrös and Miklós, 2006) a procedure for computing
223 the likelihood in a restricted gain-loss-duplication model (assuming $0 < \kappa$ and
224 $0 < \lambda < \mu$), without imposing artificial size bounds. The weakness of that pro-
225 cedure is potential numerical instability, due to the use of alternating sums in the
226 formulas. We found practical cases (such as the archaeal gene content study we
227 report below), where the numerical instability led to serious errors. The novel pro-
228 cedure presented here is numerically stable, as well as computationally efficient. It
229 applies to arbitrary gain-loss-duplication models, including degenerate cases such
230 as the one of (Hahn *et al.*, 2005) with $\lambda = \mu$ and $\kappa = 0$. The algorithm takes
231 $O(M^2n)$ time to complete for a phylogenetic profile over n species and M total
232 number of genes (see Supplemental Material).

233 **Gene content evolution in Archaea**

234 Archaea constitute one of the three main domains of cellular life, and are notable
235 for a spectacular diversity of adaptive strategies to extreme environments (Garrett
236 and Klenk, 2006). We examined gene content evolution in Archaea. For the pur-
237 poses of the study we have selected 28 completely sequenced genomes covering all
238 major physiological and metabolic groups recognized in cultured Archaea: ther-
239 mophiles, halophiles, acidophiles, nitrifiers and methanogens (Valentine, 2007).
240 Homolog gene families were extracted from the arCOG (archaeal clusters of or-
241 thologous groups) database (Makarova *et al.*, 2007), and combined with groupings
242 of genes that have no archaeal homologs outside of single genomes. The complete
243 data set consists of 14216 families, of which 7461 are among the arCOGs.

244 **Phylogenetic relationships**

245 Archaeal phylogenetic relationships have been resolved to an increasing degree of
246 confidence (Forterre *et al.*, 2006) with the aid of accumulating sequence data. Fig-
247 ure 1 shows our consensual phylogeny based on maximum likelihood trees for con-
248 catenated alignments of 46 ribosomal proteins (r-proteins) and 88 unique conserved
249 proteins (uc-proteins), which are precisely those that have exactly one homolog in
250 each sampled genome. Congruent phylogenies were proposed before (Forterre
251 *et al.*, 2006), based on complete phylogenomics evidence. In our study, r-proteins
252 and uc-proteins show solid support for most recognized phylogenetic relationships,
253 but provide contradictory signals for the placement of some euryarchaeal groups.
254 Notably, both sequence data sets support the basal position of *N. equitans*, which
255 was originally thought to be a specimen of a separate group from Euryarchaeota
256 and Crenarchaeota (Waters *et al.*, 2003), but is more likely an early-branching eu-
257 ryarchaeal organism (Makarova and Koonin, 2005; Forterre *et al.*, 2006). The data
258 also support the early branching position of non-thermophilic crenarchaea repre-
259 sented by *C. symbiosum*. In fact, non-thermophilic crenarchaea may constitute a
260 separate phylum from Euryarchaeota and Crenarchaeota, tentatively named Thau-
261 marchaeota (Brochier-Armanet *et al.*, 2008).

262 [Figure 1 about here.]

263 The observed uncertainties about euryarchaeal groups concern the placement
264 of Thermoplasmata, and so-called Class I methanogens (Baptiste *et al.*, 2005)
265 comprising Methanopyrales, Methanobacteriales and Methanococcales. Thermo-
266 plasmata were originally thought to be an early-branching lineage of Euryarchaeota
267 (Forterre *et al.*, 2006), but analyses of r-proteins (Matte-Tailliez *et al.*, 2002) have

268 provided strong evidence for their late-branching position after Class I methanogens
269 as in Figure 1. R-proteins in our study support the late-branching of Thermo-
270 plasmatales (89% bootstrap value), but a maximum-likelihood tree built from uc-
271 proteins places Thermoplasmatales between Nanoarchaea and Thermococcales (66%
272 BV). It has been argued that this placement is due to long-branch attraction (Matte-
273 Tailliez *et al.*, 2002; Brochier *et al.*, 2004), a frequent systematic bias of sequence
274 evolution models (Rodríguez-Ezpeleta *et al.*, 2007). Indeed, after we removed
275 *N. equitans* and *C. symbiosum* from the uc-protein data set, the late-branching po-
276 sition of Thermoplasmatales regained solid support (100% BV).

277 The correct phylogenetic position of *M. kandleri* (Metka) is one of the re-
278 maining puzzles in archaeal evolution. The existence of close phylogenetic re-
279 lationships between Class I methanogens is fairly certain, but different protein
280 sets and taxonomic sampling give conflicting or weak indications (Slesarev *et al.*,
281 2002; Brochier *et al.*, 2004, 2005; Gao and Gupta, 2007) about the exact branch-
282 ing order among Methanopyrales, Methanobacteriales and Methanococcales. R-
283 proteins in our study give a weak support for the monophyly of Methanococcales
284 and Methanobacteriales at the exclusion of Methanopyrales (49% BV) and faintly
285 favor the paraphyly of Class I methanogens (37% BV for the immediate split of
286 Methanopyrales between Thermococcales and Methanobacteriales/Methanococcales;
287 see Supplemental Material). Uc-proteins, however, solidly point to the monophyly
288 of Class I methanogens (> 97% BV). Interestingly, the maximum-likelihood trees
289 built from uc-proteins do not resolve well the relationships between Halobacteri-
290 ales, Methanosarcinales and Methanomicrobiales (see Supplemental Material), but
291 there is little reason to doubt that r-proteins provide a genuine phylogenetic signal
292 about the monophyly of Class II methanogens (Baptiste *et al.*, 2005; Brochier-

293 Armanet *et al.*, 2008), uniting Methanosarcinales and Methanomicrobiales.

294 We conclude that based on protein sequences, Thermoplasmatales constitute
295 a late-branching euryarchaeal lineage, and their early-branching status is a long-
296 branch attraction artifact. Furthermore, the sequences provide evidence of the
297 monophyly of both Class I and Class II methanogens.

298 **Evolutionary rates: correlations between sequence and gene content** 299 **evolution**

300 We experimented with models of increasing complexity that combine lineage- and
301 gene-specific factors in the gain-loss-duplication processes. Specifically, we as-
302 sumed that the process for family f on branch e is characterized by the rates
303 $\kappa = \hat{\kappa}_e \kappa_f$, $\mu = \hat{\mu}_e \mu_f$, $\lambda = \hat{\lambda}_e \lambda_f$, and runs for a duration of $t = \hat{t}_e t_f$. Here,
304 $\hat{t}_e, \hat{\kappa}_e, \hat{\mu}_e, \hat{\lambda}_e$ are branch-specific process parameters, and $t_f, \kappa_f, \mu_f, \lambda_f$ are family-
305 specific rate variation coefficients. Starting with simple models with invariant
306 family-specific coefficients, we introduced rate variation in a model hierarchy with
307 increasing complexity. In more complex models, some coefficients were drawn
308 randomly from a discretized Gamma distribution (Yang, 1994). Different family-
309 specific coefficients do not have the same impact on the model fit. We found the
310 largest improvement when introducing variation in edge length (t_f), followed by
311 duplication-rate variation (λ_f). Further variation in loss and gain rates led to in-
312 significant improvements in the model fit, and were not assumed in the analysis.

313 [Figure 2 about here.]

314 In the absence of extraneous scaling, we set $\hat{t}_e = 1$ in order to examine the total
315 rates of gene content change on each edge e . We found a conspicuous correlation

316 across branches between the rate of sequence evolution (expected numbers of sub-
317 stitutions per site for ribosomal proteins) and the component rates of gene content
318 evolution: on this point, see Figure 2 for loss, and the Supplemental Material for
319 duplication and gain. More precisely, the correlation holds for the lineage-specific
320 components of loss, duplication and gain rates in a decreasing order of strength
321 (P -values of $1.1 \cdot 10^{-11}$, $8.2 \cdot 10^{-6}$, $1.6 \cdot 10^{-4}$, respectively, by Student's t -test for
322 Spearman rank-order correlation coefficient).

323 The apparent correlations between gene content and sequence evolution rates
324 imply that a steady balance has been maintained between drift and natural selec-
325 tion in almost all lineages. Loss and duplication rates, in particular, have similar
326 vagaries as amino acid substitution rates, and provide thus comparable molecu-
327 lar clocks. We measured each terminal node's depth by summing the rates along
328 branches from the root to the node in question. Excluding *N. equitans* and *C. sym-*
329 *biosum*, the coefficient of variation of the depth is 26% for protein sequences, 23%
330 for gene loss rates and 20% for duplication rates. Depths by gene gain rates span
331 about a four-fold range: for substitution, loss, and duplication, the span is close to
332 two-fold.

333 Genes have thus been eliminated in all archaeal lineages with a fairly universal
334 constancy, apart from occasional accelerations. In other words, genome degrada-
335 tion processes seem to persist at a fairly common intensity in every lineage (Mira
336 *et al.*, 2001). Conceivably, genome decay is counterbalanced by natural selection
337 that eliminates deleterious mutations. The root cause of dramatically increased
338 gene loss in endosymbionts such as *N. equitans* (Makarova and Koonin, 2005) may
339 be reduced selection (Hershberg *et al.*, 2007; Koonin and Wolf, 2008). Principles
340 of population genetics imply that changes in population size alone can explain rate

341 changes (Lynch, 2006): selection power is weaker in a smaller population, which
342 should manifest in accelerated evolution of sequences (Ohta, 1972) and gene con-
343 tent.

344 We examined the differences between evolutionary rates in sibling terminal
345 taxa for signs of natural selection. Figure 2 shows that gene loss and amino
346 acid substitution rates differ in a concerted fashion for three pairs, that is, for
347 *M. stadtmanæ*-*M. thermoautotrophicus*, *Halobacterium* sp.-*H. marismortuimi*, and
348 *S. acidocaldarius*-*S. sulfolobus*. In seven other pairs, loss rates are essentially the
349 same, even if substitution rates may differ. The agreements between substitution
350 and gene loss rate changes attest to common selection forces and mutation pro-
351 cesses acting on different forms of genome decay, and are predicted by population-
352 genetic arguments (Lynch, 2006).

353 In the lineage leading to *M. stadtmanæ*, a human commensal (Fricke *et al.*,
354 2006), all rates are simultaneously larger when compared to its sibling lineage
355 *M. thermoautotrophicus*, which may be attributed to a smaller population size for
356 the former, which has a smaller habitat. Gene gain and duplication rates behave
357 in general less predictably: numerical differences between loss, gain, and duplica-
358 tion rates on sibling lineages occur in almost all possible sign combinations. The
359 observed fluctuations corroborate the intuition that selection pressures acting on
360 gain and duplication are strong and variable (Wolf *et al.*, 2002). It is plausible that
361 during episodes of massive adaptation, the selective advantages of gene acquisition
362 may outweigh possible negative consequences of an increased genome, and thus
363 drive elevated gene gains, especially if coupled with small population sizes. In our
364 case, unusually large gain rates are inferred on some of the deepest branches (such
365 as the one leading to node E1 on Figure 1 or to the halobacterial ancestor), as well

366 as on the terminal branches leading to *M. acetivorans* (Metac), *H. marismortuimi*
367 (Halma) and *P. aerophilum* (Pyrae).

368 **History of archaeal gene census: streamlining and surges**

369 We inferred a probable history of archaeal gene content using posterior probabili-
370 ties for ancestral family sizes and family size changes, computed from the phylo-
371 genetic profiles in the fitted model. Figure 3 summarizes the results by lineages.
372 (See Supplemental Material for bootstrap confidence intervals: the uncertainty in
373 ancestral family counts is estimated to be within $\pm 19\%$ for all nodes.)

374 [Figure 3 about here.]

375 Our reconstruction suggests a recurrent theme in archaeal evolution: a major
376 physiological or metabolic invention leads to a successful founding population in
377 a new environment, which then further diversifies by genomic streamlining. We
378 can see notably that Figure 3 shows only a few branches where gains prevail over
379 losses (i.e., at least twice as many gains as losses): such is the case for some deep
380 crenarchaeal and euryarchaeal branches, and the terminal lineages for *M. acetivo-*
381 *rans* and *H. marismortuimi*. About half of the remaining terminal lineages and
382 two-thirds of remaining deep lineages are dominated by loss. Moreover, there is
383 only one ancestral node (the crenarchaeal ancestor) in the entire tree for which gain
384 is dominant in both descendant lineages.

385 Why would gene loss be so prevalent? We speculate that the versatility of
386 a large genome in such extant lineages as *M. acetivorans* (Galagan *et al.*, 2002)
387 and *H. marismortuimi* (Baliga *et al.*, 2004) can be upheld for only relatively short
388 time periods. Genetic drift already leads to the diversification of descendant lin-

389 eages, which are frequently isolated, given the disconnectedness of the extreme
390 environments they dwell in (Whitaker *et al.*, 2003; Escobar-Páramo *et al.*, 2005).
391 Specialization and the loss of dispensable functions should be favorable in the de-
392 scendants that are typically under significant energy stress (Valentine, 2007). Ge-
393 nomic streamlining should also be favored by population-size effects due to the
394 isolation (Lynch, 2006), even in the case of slightly deleterious loss of function.

395 After the crenarchaeal split, the main euryarchaeal lineage has been charac-
396 terized by the accumulation of new families, culminating in a large surge on the
397 branch leading to node E1, where many new families appeared. The time interval
398 (judging by sequence divergence in Figure 1) and the extent of gene gain is similar
399 to what is seen with *H. marismortuimi* (Halma) and *M. acetivorans* (Metac). The
400 inference of large gains in the E1 lineage is due to the large number of gene fami-
401 lies shared between multiple descendant lineages, and especially between the two
402 classes of methanogens (Slesarev *et al.*, 2002; Bapteste *et al.*, 2005; Gao and Gupta,
403 2007; Makarova *et al.*, 2007). In fact, this lineage may very well have been where
404 hydrogenotrophic methanogenesis was invented, which then underwent modifica-
405 tions, extensions and degradations in subsequent lineages. It was noted in previous
406 genome-scale comparisons (Bapteste *et al.*, 2005; Gao and Gupta, 2007) that it is
407 likely that euryarchaeal lineages acquired methanogenesis predominantly by verti-
408 cal inheritance, because the associated pathways are fairly complex, and neither the
409 sequences nor the phylogenetic profiles show evidence of substantial amounts of
410 lateral gene transfer. Figure 3 suggests that methanogenesis appeared after the split
411 of Thermococcales in the company of more than 760 genes. Based on extant ex-
412 amples of archaea with such swelled genomes (Galagan *et al.*, 2002; Baliga *et al.*,
413 2004), it is plausible that the corresponding archaeal organisms were extremely

414 versatile.

415 Our inference of ancestral gene content is quite different from previous recon-
416 structions based on parsimony principles (Makarova *et al.*, 2007; Csűrös, 2008): at
417 deep nodes, we postulate larger genomes. Parsimonious reconstructions (Mirkin
418 *et al.*, 2003; Kunitz *et al.*, 2005; Csűrös, 2008) aim to minimize the number of im-
419 plied loss and gain events. As a consequence, parsimony inherently underestimates
420 the age of gene families. A probabilistic framework, such as a phylogenetic birth-
421 and-death model, makes it feasible to take all possibilities into consideration in a
422 mathematically sound way. A case in point is the last archaeal common ancestor
423 (LACA), where only about 1300 families are inferred to have been present with a
424 posterior probability of at least 90%, which is close to a parsimony-based infer-
425 ence of about 1000 families (Makarova *et al.*, 2007). Given the uncertainties of
426 most family histories, the exact genome composition of LACA is hard to estimate,
427 but the fractional probabilities point to a genome with slightly more than 2000
428 families, which is similar to such extant organisms as *S. sulfolobus*. Such a large
429 genome size implies that LACA's genomic complexity was even greater than pre-
430 viously imagined (Makarova *et al.*, 2007), on a par with modern, moderately-sized
431 archaeal genomes.

432 **Acknowledgments**

433 This work has been supported by a grant from the Natural Sciences and Engineer-
434 ing Research Council of Canada. Part of the study was done while M.Cs. was a
435 sabbatical visitor at the Rényi Institute of Mathematics, supported by a Marie-Curie
436 Transfer-of-Knowledge fellowship. We are grateful to Yuri Wolf for providing data

437 on lineage-specific gene families. We thank Igor Rogozin, Csaba Pál and Balázs
438 Papp for informative discussions.

439 **References**

- 440 Alexeyenko, A., Tamas, I., Liu, G., and Sonnhammer, E. L. L. (2006). Automatic
441 clustering of orthologs and inparalogs shared by multiple genomes. *Bioinfor-*
442 *matics*, **22**, e9–e15.
- 443 Baliga, N. S. *et al.* (2004). Genome sequence of *Haloarcula morismurtoimi*: a
444 halophilic archaeon from the Dead Sea. *Genome Research*, **14**, 2221–2234.
- 445 Bapteste, É., Brochier, C., and Boucher, Y. (2005). Higher-level classification of
446 the Archaea: evolution of methanogenesis and methanogens. *Archaea*, **1**, 353–
447 363.
- 448 Boucher, Y., Douady, C. J., Papke, R. T., Walsh, D. A., Boudreau, M. E. R., Nesbo,
449 C. L., Case, R. J., and Doolittle, W. F. (2003). Lateral gene transfer and the
450 origin of prokaryotic groups. *Annual Review of Genetics*, **37**, 283–328.
- 451 Brochier, C., Forterre, P., and Gribaldo, S. (2004). Archaeal phylogeny based
452 on proteins of the transcription and translation machineries: tackling the
453 *Methanopyrus paradox*. *Genome Biology*, **5**, R17.
- 454 Brochier, C., Forterre, P., and Gribaldo, S. (2005). An emerging phylogenetic core
455 of Archaea: phylogenies of transcription and translation machineries converge
456 following addition of new genome sequences. *BMC Evolutionary Biology*, **5**,
457 36.
- 458 Brochier-Armanet, C., Boussau, B., Gribaldo, S., and Forterre, P. (2008).
459 Mesophilic crenarchaeota: proposal for a third archaeal phylum, the Thaumarchaeota.
460 *Nature Reviews Microbiology*, **6**, 245–252.

- 461 Csűrös, M. (2008). Ancestral reconstruction by asymmetric Wagner parsimony
462 over continuous characters and squared parsimony over distributions. *Springer*
463 *Lecture Notes in Bioinformatics*, **5267**, 72–86. Proc. Sixth RECOMB Compar-
464 ative Genomics Satellite Workshop.
- 465 Csűrös, M. and Miklós, I. (2006). A probabilistic model for gene content evolu-
466 tion with duplication, loss, and horizontal transfer. *Springer Lecture Notes in*
467 *Bioinformatics*, **3909**, 206–220. Proc. Tenth Annual International Conference
468 on Research in Computational Molecular Biology (RECOMB).
- 469 Edgar, R. C. (2004). MUSCLE: multiple sequence alignment with high accuracy
470 and high throughput. *Nucleic Acids Research*, **32**(5), 1792–1797.
- 471 Escobar-Páramo, P., Gosh, S., and DiRuggiero, J. (2005). Evidence for genetic
472 drift in the diversification of a geographically isolated population of the hyper-
473 thermophylic archaeon *Pyrococcus*. *Molecular Biology and Evolution*, **22**(11),
474 2297–2303.
- 475 Felsenstein, J. (1981). Evolutionary trees from DNA sequences: A maximum like-
476 lihood approach. *Journal of Molecular Evolution*, **17**, 368–376.
- 477 Felsenstein, J. (1992). Phylogenies from restriction sites, a maximum likelihood
478 approach. *Evolution*, **46**, 159–173.
- 479 Fitz-Gibbon, S. T. and House, C. H. (1999). Whole genome-based phylogenetic
480 analysis of free-living microorganisms. *Nucleic Acids Research*, **27**(21), 4218–
481 4222.

- 482 Forterre, P., Gribaldo, S., and Brochier-Armanet, C. (2006). Natural history of the
483 archaeal domain. In Garrett and Klenk (2006), chapter 2, pages 17–28.
- 484 Fricke, W. F. *et al.* (2006). The genome sequence of *Methanosphaera stadtmanae*
485 reveals why this human intestinal archaeon is restricted to methanol and H₂ for
486 methane formation and ATP synthesis. *Journal of Bacteriology*, **188**(2), 642–
487 658.
- 488 Galagan, J. E., Nusbaum, C., Roy, A., *et al.* (2002). The genome of *M. acetivorans*
489 reveals extensive metabolological and physiological diversity. *Genome Research*,
490 **12**, 532–542.
- 491 Gao, B. and Gupta, R. S. (2007). Phylogenomic analysis of proteins that are dis-
492 tinctive of Archaea and its main subgroups and the origin of methanogenesis.
493 *BMC Genomics*, **8**, 86.
- 494 Garrett, R. A. and Klenk, H.-P., editors (2006). *Archaea: Evolution, Physiology,*
495 *and Molecular Biology*. Blackwell Publishing, Malden, Mass.
- 496 Gogarten, J. P. and Townsend, J. P. (2005). Horizontal gene transfer, genome inno-
497 vation and evolution. *Nature Reviews Microbiology*, **3**, 679–687.
- 498 Guindon, S. and Gascuel, O. (2003). A simple, fast, and accurate accurate algo-
499 rithm to estimate large phylogenies by maximum likelihood. *Systematic Biology*,
500 **52**(5), 696–704.
- 501 Hahn, M. W., De Bie, T., Stajich, J. E., Nguyen, C., and Cristianini, N. (2005).
502 Estimating the tempo and mode of gene family evolution from comparative ge-
503 nomic data. *Genome Research*, **15**, 1153–1160.

- 504 Henikoff, S., Greene, E. A., Pietrokovski, S., Bork, P., Atwood, T. K., and Hood, L.
505 (1997). Gene families: the taxonomy of protein paralogs and chimeras. *Science*,
506 **278**, 609–614.
- 507 Hershberg, R., Tang, H., and Petrov, D. A. (2007). Reduced selection leads to
508 accelerated gene loss in *Shigella*. *Genome Biology*, **8**, R164.
- 509 Iwasaki, W. and Takagi, T. (2007). Reconstruction of highly heterogeneous gene-
510 content evolution across the three domains of life. *Bioinformatics*, **23**(13), i230–
511 i239.
- 512 Jordan, M. I. (2004). Graphical models. *Statistical Science*, **19**(1), 140–155.
- 513 Karlin, S. and McGregor, J. (1958). Linear growth, birth, and death processes.
514 *Journal of Mathematics and Mechanics*, **7**(4), 643–662.
- 515 Kendall, D. G. (1949). Stochastic processes and population growth. *Journal of the*
516 *Royal Statistical Society Series B*, **11**(2), 230–282.
- 517 Koonin, E. V. (2005). Orthologs, paralogs, and evolutionary genomics. *Annual*
518 *Review of Genetics*, **39**, 309–338.
- 519 Koonin, E. V. and Wolf, Y. I. (2008). Genomics of bacteria and archaea: the emerg-
520 ing dynamic view of the prokaryotic world. *Nucleic Acids Research*. Advance
521 access published online on October 23, 2008. DOI:10.1093/nar/gkn668.
- 522 Kunin, V., Goldovsky, L., Darzentas, N., and Ouzounis, C. A. (2005). The net
523 of life: reconstructing the microbial phylogenetic network. *Genome Research*,
524 **15**(7), 954–959.

- 525 Lynch, M. (2006). Streamlining and simplification of microbial genome achitec-
526 ture. *Annual Review of Microbiology*, **60**, 327–349.
- 527 Makarova, K. and Koonin, E. V. (2005). Evolutionary and functional genomics of
528 the Archaea. *Current Opinion in Microbiology*, **8**, 586–594.
- 529 Makarova, K. S., Sorokin, A. V., Novichkov, P. S., Wolf, Y. I., and Koonin, E. V.
530 (2007). Clusters of orthologous genes for 41 archaeal genomes and implications
531 for evolutionary genomics of archaea. *Biology Direct*, **2**, 33.
- 532 Matte-Tailliez, O., Brochier, C., Forterre, P., and Philippe, H. (2002). Archaeal
533 phylogeny based on ribosomal proteins. *Molecular Biology and Evolution*,
534 **19**(5), 631–639.
- 535 Mira, A., Ochman, H., and Moran, N. A. (2001). Deletional bias and the evolution
536 of bacterial genomes. *Trends in Genetics*, **17**(10), 589–596.
- 537 Mirkin, B. G., Fenner, T. I., Galperin, M. Y., and Koonin, E. V. (2003). Algorithms
538 for computing evolutionary scenarios for genome evolution, the last universal
539 common ancestor and dominance of horizontal gene transfer in the evolution of
540 prokaryotes. *BMC Evolutionary Biology*, **3**, 2.
- 541 Nei, M. and Rooney, A. P. (2005). Concerted and birth-and-death evolution of
542 multigene families. *Annual Review of Genetics*, **39**(1), 121–152.
- 543 Ohta, T. (1972). Population size and rate of evolution. *Journal of Molecular*
544 *Evolution*, **1**, 305–314.
- 545 Pellegrini, M., Marcotte, E. M., Thompson, M. J., Eisenberg, D., and Yeates, T. O.
546 (1999). Assigning protein functions by comparative genome analysis: protein

547 phylogenetic profiles. *Proceedings of the National Academy of Sciences of the*
548 *USA*, **96**(8), 4285–4288.

549 Press, W. H., Teukolsky, S. A., Vetterling, W. V., and Flannery, B. P. (1997). *Nu-*
550 *merical Recipes in C: The Art of Scientific Computing*. Cambridge University
551 Press, second edition.

552 Rodríguez-Ezpeleta, N., Brinkmann, H., Roure, B., Lartillot, N., Lang, B. F., and
553 Philippe, H. (2007). Detecting and overcoming systematic errors in genome-
554 scale phylogenies. *Systematic Biology*, **56**(3), 389–399.

555 Ross, S. M. (1996). *Stochastic Processes*. Wiley & Sons, second edition.

556 Slesarev, A. I. *et al.* (2002). The complete genome of hyperthermophile
557 *Methanopyrus kandleri* AV19 and monophyly of archaeal methanogens. *Pro-*
558 *ceedings of the National Academy of Sciences of the USA*, **99**(7), 4644–4649.

559 Snel, B., Bork, P., and Huynen, M. A. (1999). Genome phylogeny based on gene
560 content. *Nature Genetics*, **21**(1), 108–110.

561 Snel, B., Bork, P., and Huynen, M. A. (2002). Genomes in flux: the evolution of
562 archaeal and proteobacterial gene content. *Genome Research*, **12**(1), 17–25.

563 Sonnhammer, E. L. L. and Koonin, E. V. (2002). Orthology, paralogy and proposed
564 classification for paralog subtypes. *Trends in Genetics*, **18**(12), 619–620.

565 Spencer, M., Susko, E., and Roger, A. J. (2006). Modelling prokaryote gene con-
566 tent. *Evolutionary Bioinformatics Online*, **2**, 165–186.

567 Takács, L. (1962). *Introduction to the Theory of Queues*. Oxford University Press,
568 New York.

- 569 Tatusov, R. L., Koonin, E. V., and Lipman, D. J. (1997). A genomic perspective on
570 protein families. *Science*, **278**, 631–637.
- 571 Tekaiia, F., Lazcano, A., and Dujon, B. (1999). The genomic tree as revealed from
572 whole proteome comparisons. *Genome Research*, **9**(6), 550–557.
- 573 Valentine, D. L. (2007). Adaptations to energy stress dictate the ecology and evo-
574 lution of the Archaea. *Nature Reviews Microbiology*, **5**, 316–323.
- 575 Waters, E. *et al.* (2003). The genome of Nanoarchaeum equitans: insights into
576 early archaeal evolution and derived parasitism. *Proceedings of the National*
577 *Academy of Sciences of the USA*, **100**(22), 12984–12988.
- 578 Whitaker, R. J., Grogan, D. W., and Taylor, J. W. (2003). Geographic barriers
579 isolate endemic populations of hyperthermophilic archaea. *Science*, **301**(5635),
580 976–978.
- 581 Wolf, Y. I., Rogozin, I. B., Grishin, N. V., and Koonin, E. V. (2002). Genome trees
582 and the Tree of Life. *Trends in Genetics*, **18**(9), 472–479.
- 583 Yang, Z. (1994). Maximum likelihood phylogenetic estimation from DNA se-
584 quences with variable rates over sites: approximate methods. *Journal of Molec-*
585 *ular Evolution*, **39**, 306–314.

586 **List of Figures**

587 1 Consensus evolutionary tree of Archaea in the study. The consen-
588 sus is based on maximum-likelihood trees for concatenated align-
589 ments of ribosomal and unique conserved proteins. Branch lengths
590 are set by maximum likelihood for the 88 unique conserved pro-
591 teins. Recognized archaeal orders are highlighted. The boxed
592 triples on the left show the percentage of bootstrap samples sup-
593 porting the particular edges in three data sets (from 500 replicates
594 for each set): r-proteins, uc-proteins, and uc-proteins without *C. sym-*
595 *biosum* (**Censy**) and *N. equitans* (**Naneq**). All other edges have
596 > 97% bootstrap support in all data sets. 30

597 2 Branch-specific loss rates $\hat{\mu}_e \hat{t}_e$ compared to expected numbers of
598 substitutions (or *edge length*) for each branch *e*. Pairs of sibling
599 terminal taxa are connected by lines. 31

600 3 A digest of gene content evolution in Archaea. The bar graphs
601 plot posterior means for number of families. The chart **on the left**
602 shows the number of families with at least one homolog; the fatter
603 part of the bar is proportional to the number of multi-gene families.
604 The chart **in the middle** plots the families acquired and lost on the
605 branch leading to the indicated node. The net change is highlighted
606 by the solid part of the bars. The chart **on the right** shows how
607 many families underwent a contraction from multi-gene to single-
608 gene composition, or expanded from a single homolog to multiple
609 paralogs. Note that scaling is the same on the left-hand side and in
610 the middle, but different on the right-hand side. 32

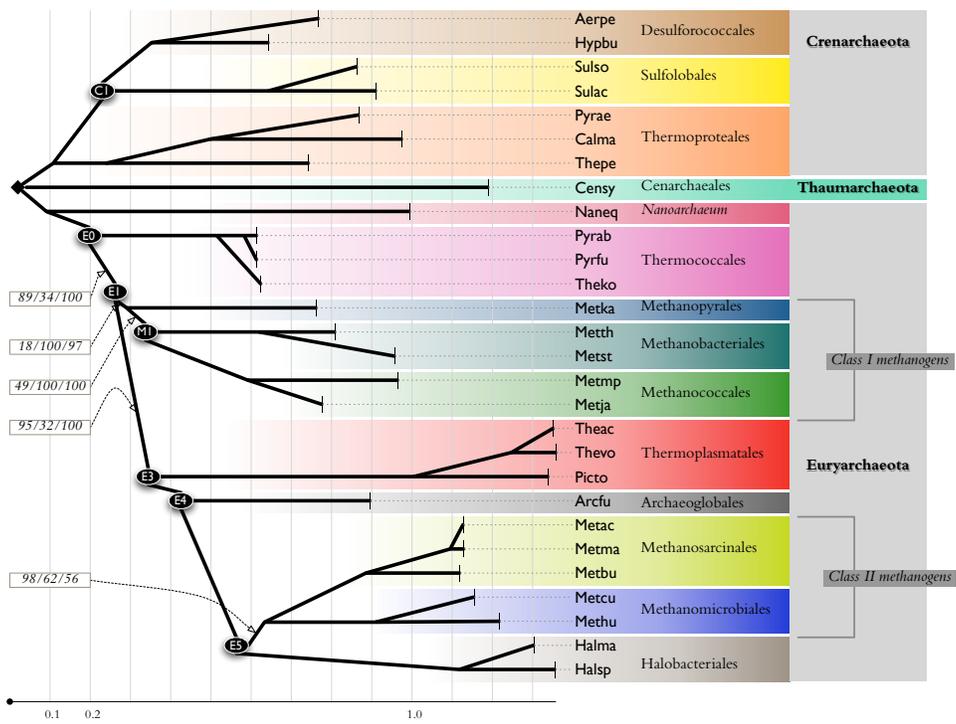


Figure 1

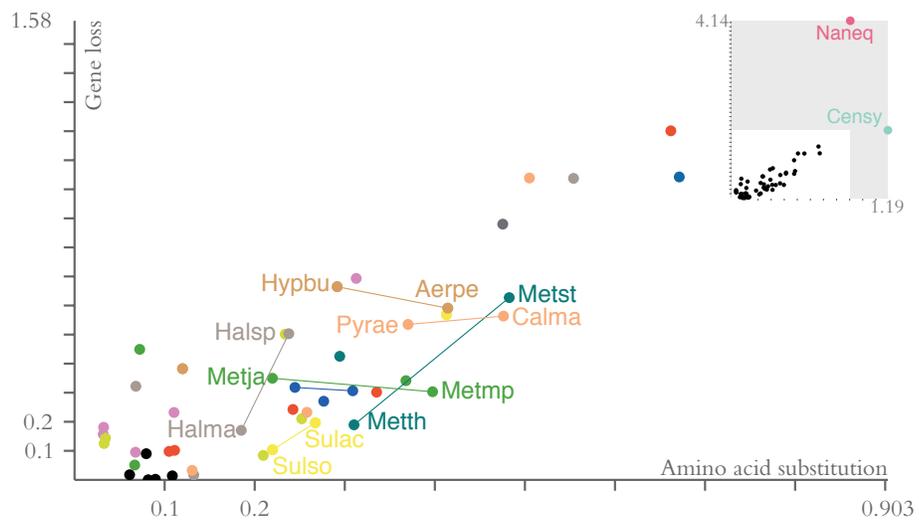


Figure 2

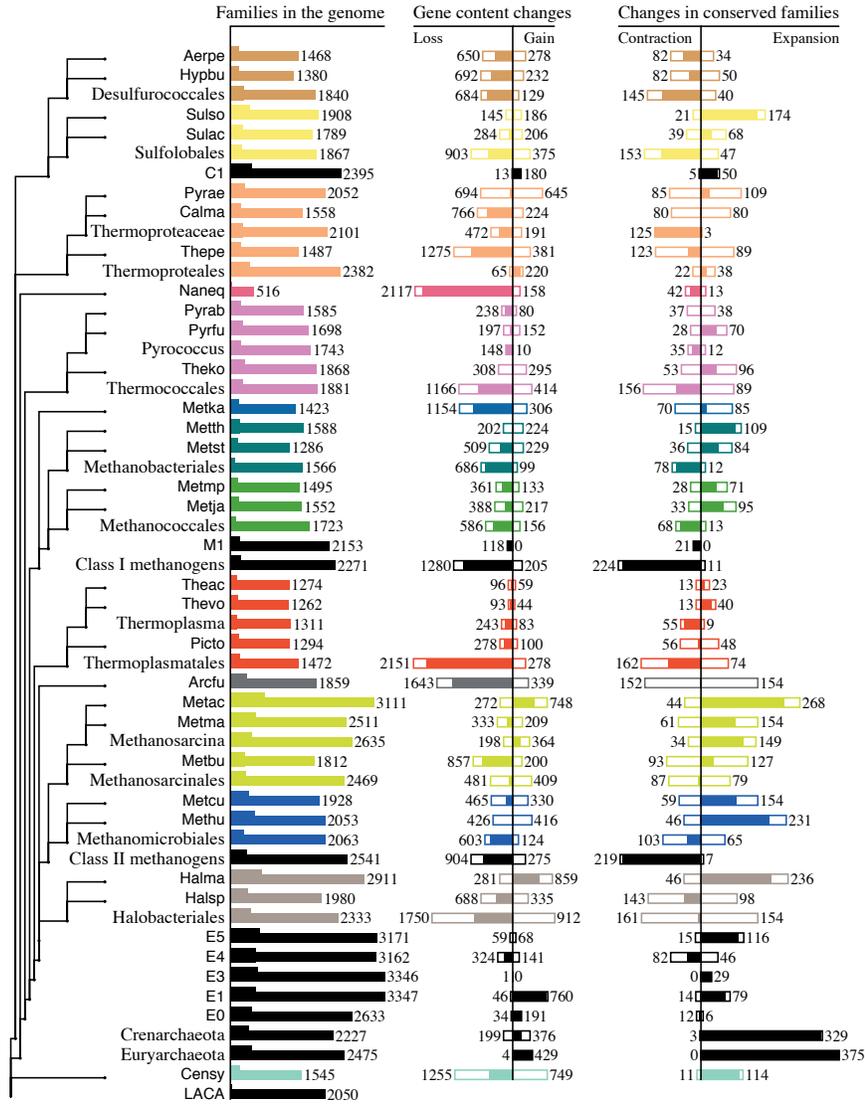


Figure 3