

1 Order of amino acid recruitment into the genetic code resolved by
2 Last Universal Common Ancestor's protein domains

3 Sawsan Wehbi¹, Andrew Wheeler¹, Benoit Morel², Nandini Manepalli³, Bui Quang Minh⁴, Dante
4 S. Laretta⁵, Joanna Masel⁶

5 ¹Genetics Graduate Interdisciplinary Program, University of Arizona, Tucson, Arizona, 85721,
6 USA

7 ²Computational Molecular Evolution Group, Heidelberg Institute for Theoretical Studies,
8 Heidelberg, Germany

9 ³Department of Molecular and Cellular Biology, University of Arizona, Tucson, AZ, 85721, USA

10 ⁴School of Computing, Australian National University, Canberra, ACT, Australia

11 ⁵Lunar and Planetary Laboratory, University of Arizona, Tucson, AZ 85721, USA

12 ⁶Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, AZ, 85721,
13 USA

14 **Corresponding author:** Joanna Masel

15 **Email:** masel@arizona.edu

16 **Author Contributions:** S.W. and J.M. conceived the study, with input from D.S.L. N.M. produced
17 Figure 3, and S.W. conducted all other analyses, with advice from J.M., A.W., B.Q.M., and B.M.
18 B.M. implemented a new feature in GeneRax to assist the analyses. S.W. wrote the first draft of
19 the manuscript, with subsequent editing from S.W., J.M., D.S.L., A.W., and B.Q.M.

20 **Competing Interest Statement:** The authors declare no competing interest.

21 **Classification:** Biological Sciences, Evolution.

22 **Keywords:** Origins of life, early life, astrobiology, phylostratigraphy, antioxidants, metalloproteins

23 **Abstract**

24 The current “consensus” order in which amino acids were added to the genetic code is based on
25 potentially biased criteria, such as absence of sulfur-containing amino acids from the Urey-Miller
26 experiment which lacked sulfur. More broadly, abiotic abundance might not reflect biotic
27 abundance in the organisms in which the genetic code evolved. Here, we instead identify which
28 protein domains date to the last universal common ancestor (LUCA), then infer the order of
29 recruitment from deviations of their ancestrally reconstructed amino acid frequencies from the
30 still-ancient post-LUCA controls. We find that smaller amino acids were added to the code earlier,
31 with no additional predictive power in the previous “consensus” order. Metal-binding (cysteine and
32 histidine) and sulfur-containing (cysteine and methionine) amino acids were added to the genetic
33 code much earlier than previously thought. Methionine and histidine were added to the code
34 earlier than expected from their molecular weights, and glutamine later. Early methionine
35 availability is compatible with inferred early use of S-adenosylmethionine, and early histidine with
36 its purine-like structure and the demand for metal-binding. Even more ancient protein sequences
37 — those that had already diversified into multiple distinct copies prior to LUCA — have
38 significantly higher frequencies of aromatic amino acids (tryptophan, tyrosine, phenylalanine and
39 histidine), and lower frequencies of valine and glutamic acid than single copy LUCA sequences. If
40 at least some of these sequences predate the current code, then their distinct enrichment
41 patterns provide hints about earlier, alternative genetic codes.

42 **Significance Statement**

43 The order in which the amino acids were added to the genetic code was previously inferred from
44 consensus among forty metrics. Many of these reflect abiotic abundance on ancient Earth.
45 However, the abundances that matter are those within primitive cells that already had
46 sophisticated RNA and perhaps peptide metabolism. Here, we directly infer the order of
47 recruitment from the relative ancestral amino acid frequencies of ancient protein sequences.
48 Small size predicts ancient amino acid enrichment better than the previous consensus metric
49 does. We place metal-binding and sulfur-containing amino acids earlier than previously thought,
50 highlighting the importance of metal-dependent catalysis and sulfur metabolism to ancient life.
51 Understanding early life has implications for our search for life elsewhere in the universe.

52

53

54 **Main Text**

55 **Introduction**

56 The modern genetic code was likely assembled in stages, hypothesized to begin with “early”
57 amino acids present on Earth before the emergence of life (possibly delivered by extraterrestrial
58 sources such as asteroids or comets), and ending with “late” amino acids requiring biotic
59 synthesis (1, 2). For example, the Urey-Miller experiment (3) has been used to identify which
60 amino acids were available abiotically and are thus likely to have come earlier than those
61 requiring biotic synthesis. The order of amino acid recruitment, from early to late, was inferred by
62 taking statistical consensus among 40 different rankings (4), none of which constitute strong
63 evidence on their own. On the basis of this ordering, Moosmann (5) hypothesized that the first
64 amino acids recruited into the genetic code were those that were useful for membrane anchoring,
65 then those useful for halophilic folding, then for mesophilic folding, then for metal binding, and
66 finally for their antioxidant properties. However, a late role for metal-binding amino acids is
67 puzzling; many metalloproteins date back to the Last Universal Common Ancestor’s (LUCA)’s
68 proteome, where they are presumed to be key to the emergence of biological catalysis (6).

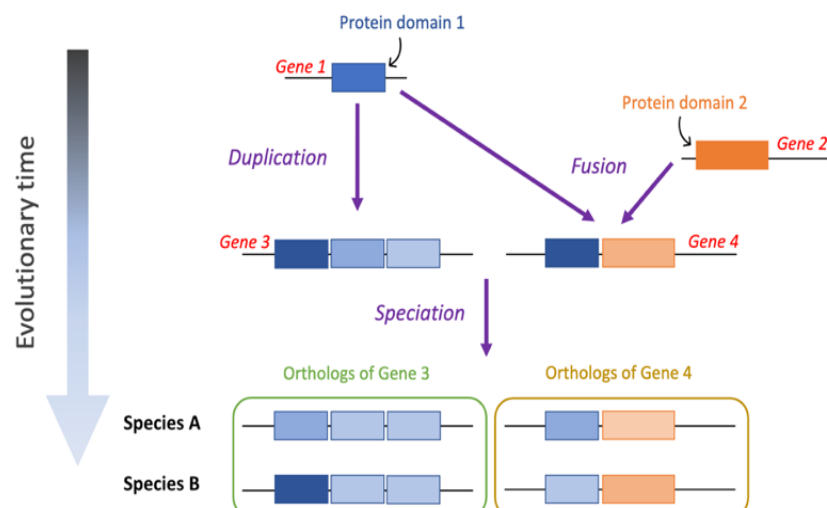
69 Indeed, the late status of some amino acids is disputed (7). For example, the Urey-Miller
70 experiment (3) did not include sulfur, and so should not have been used to infer that the sulfur-
71 containing amino acids cysteine and methionine were late additions. Methionine and
72 homocysteine (a product of cysteine degradation) were detected in hydrogen sulfide (H₂S)-rich

73 spark discharge experiments, suggesting that methionine and cysteine could be abiotically
74 produced (8). A nitrile-activated dehydroalanine pathway can produce cysteine from abiotic serine
75 that is produced from a Strecker reaction (9), further demonstrating the possibility of its early
76 chemical availability.

77 Histidine's classification as abiotically unavailable also contributed to its annotation as late (4).
78 While histidine can be abiotically synthesized from erythrose reacting with formamidine followed
79 by a Strecker synthesis reaction (10), the reactant concentrations might have been insufficient in
80 a primitive earth environment (11). More importantly, because histidine resembles a purine, even
81 if histidine were abiotically unavailable, it might have had cellular availability at the time of genetic
82 code construction (12), in an organism that biotically synthesized ribosomes, and that might also
83 have already utilized amino acids and peptides. Indeed, histidine is the most commonly
84 conserved residue in the active site of enzymes (13).

85 To directly infer the order of recruitment from protein sequence data, without reference to abiotic
86 availability arguments, we consider that some of LUCA's proteins were born prior to the
87 completion of the genetic code (14). We predict that ancestrally reconstructed sequences from
88 this era will be enriched in early amino acids and depleted in late amino acids. Previous analyses
89 relied on conserved residues within a small number of LUCA proteins (15, 16). Here, we classify
90 a larger set of protein-coding domains that date back to LUCA, rather than being more recently
91 born, e.g., *de novo* from non-coding sequences or alternative reading frames (17, 18). We
92 compare reconstructed ancient amino acid frequencies of the most ancient vs. moderately
93 ancient protein cohorts, to deduce the order in which amino acids were incorporated into the
94 genetic code.

95 We take advantage of gene-tree species-tree reconciliation methods (19) to infer LUCA's protein
96 sequences. Previous analyses focused on the age of orthologous gene families (20-22); ours is
97 the first to infer which protein domains date back to LUCA. Protein domains are the basic units of
98 proteins, that can fold, function, and evolve independently (23). Proteins often contain multiple
99 protein domains, each of which might have a different age (Figure 1). For the purpose of inferring
100 ancient amino acid usage, what matters is the age of the protein domain, not that of the whole
101 protein that it is part of. We use protein domain annotations from the Pfam database (24). We
102 recognize Pfams present in LUCA by trimming horizontal gene transfer (HGT) events, and by
103 exploiting long archaeal-bacterial branches (Figure 2; see Methods for details).

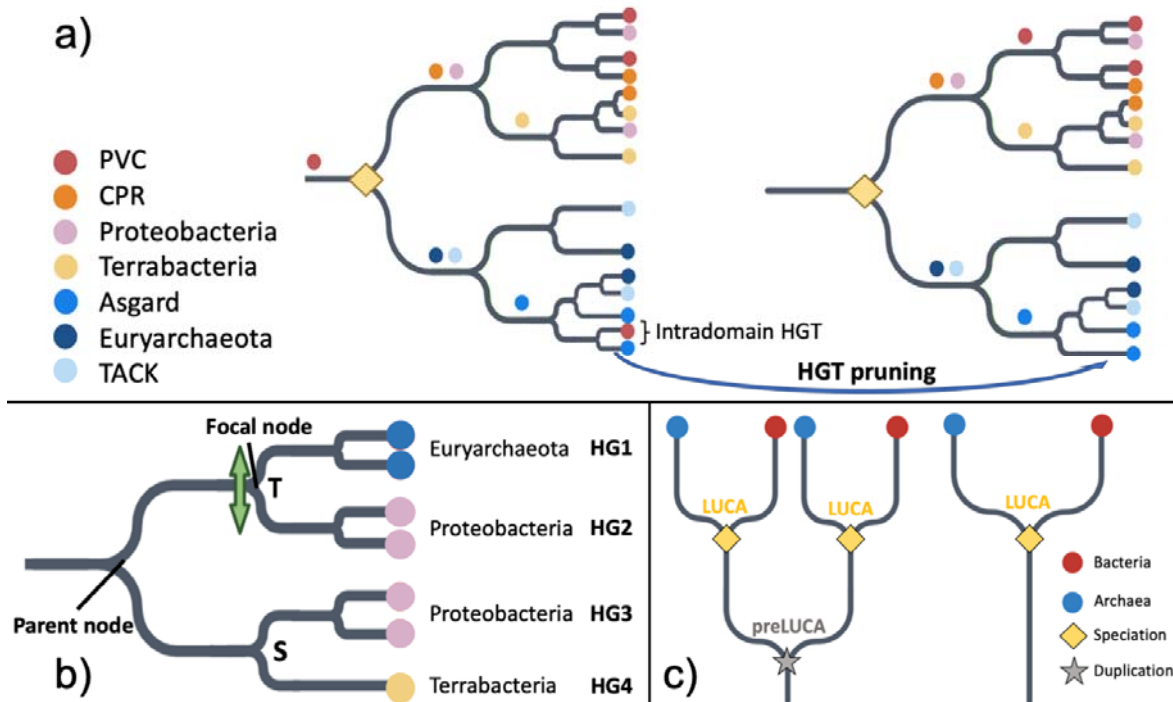


104

105 **Figure 1. The evolutionary history of a protein domain may date back further in time than**
106 **that of the whole-gene ortholog that it is part of. Multi-domain genes 3 and 4 originated**

107 around the same time. However, they are made up of two protein domains (blue & orange boxes)
 108 that emerged and diverged at different points in time – domain 1 is older than domain 2.

109



110

111

112 **Figure 2. Criteria for (a) LUCA Pfam annotation, (b) identifying HGT to be filtered, and (c)**

113 **pre-LUCA Pfam annotation.** Details are in Methods, with a brief summary here. a) Pruning HGT

114 between archaea and bacteria reveals a LUCA node as dividing bacteria and archaea at the root.

115 Colored circles are indicated just upstream of the most recent common ancestor (MRCA) of all

116 copies of that Pfam found within the same taxonomic supergroup. We recognize a total of five

117 bacterial supergroups (FCB, PVC, CPR, Terrabacteria and Proteobacteria (75, 76)) and four

118 archaeal supergroups (TACK, DPANN, Asgard and Euryarchaeota (77, 78)); only 4 out of 5

119 bacterial supergroups and 3 out of 4 archaeal supergroups are shown. The yellow diamond

120 indicates LUCA as a speciation event between archaea and bacteria. We do not assume that the

121 LUCA coalescence timing was the same for every Pfam (94). Prior to HGT pruning, PVC

122 sequences can be found on either side of the two lineages divided by the root. After pruning

123 intradomain HGT, four MRCAs are found one node away from the root, and three more MRCAs

124 are found two nodes away from the root, fulfilling our other LUCA criterion described in the

125 Methods, namely presence of at least three bacterial and at least two archaeal supergroup

126 MRCAs one to two nodes away from the root. b) Criteria for pruning likely HGT between archaea

127 and bacteria (see Methods for details). We partition into monophyletic groups of sequences in the

128 same supergroup; in this example, there are four such groups, representing two bacterial

129 supergroups and one archaeal supergroup. There is one 'mixed' node, separating an archaeal

130 group (HG1) from a bacterial group (HG2). It is also annotated by GeneRax (19) as a transfer 'T'.

131 The bacterial nature of groups 3 and 4 indicates a putative HGT direction from group 2 to group

132 1. Group 2 does not contain any Euryarchaeota sequences, meeting the third and final

133 requirement for pruning of group 1. If neither Proteobacteria or Euryarchaeota sequences were

134 present among the other descendants of the parent node, both groups 1 and 2 would be

135 considered acceptors of a transferred Pfam and would both be pruned from the tree. c) Pre-LUCA

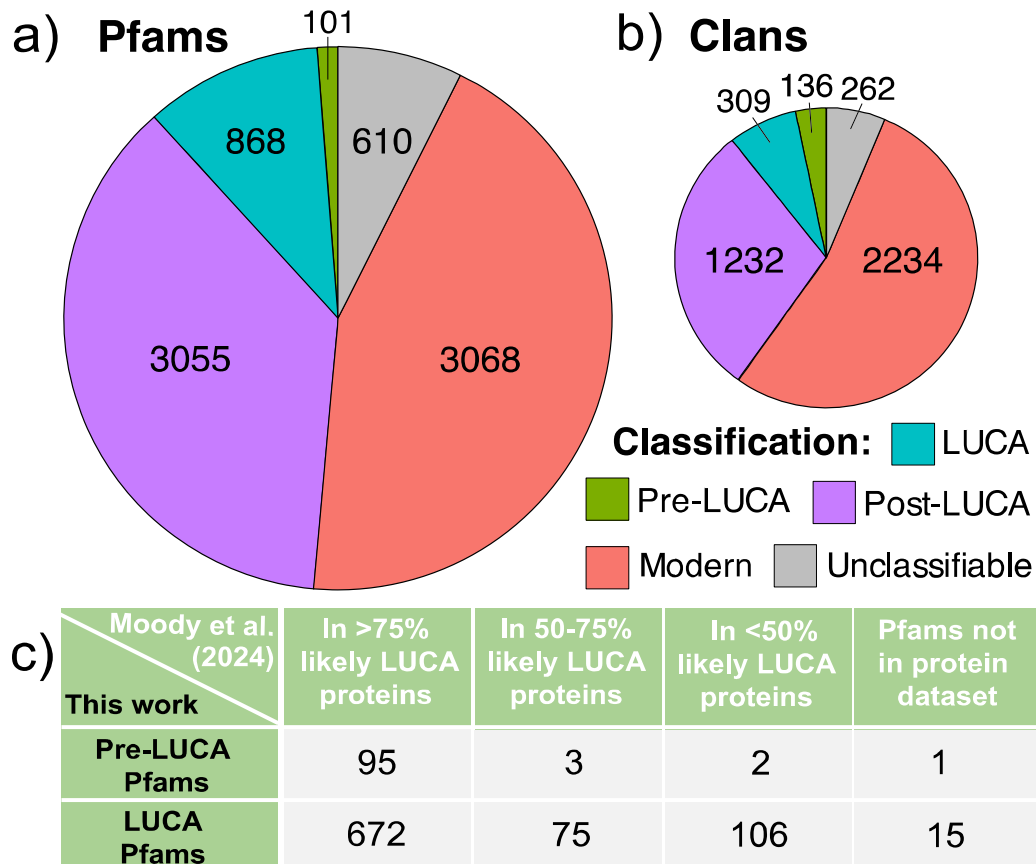
136 Pfams have at least two nodes annotated as LUCA.

137

Results

138 **Ancient protein domain classifications agree with whole-gene classifications**

139 We classify 969 Pfams and 445 clans (sets of one or more Pfams that are evolutionary related)
 140 as present in LUCA (Figures 3a and 3b; detailed lists in Supplementary Tables 1 and 2). We
 141 compare these to the 3055 Pfams and 1232 clans that we classify as ancient but post-LUCA
 142 (including Last Bacterial Common Ancestor (LBCA) and Last Archaeal Common Ancestor (LACA)
 143 candidates). Encouragingly, 88.6% of Pfams that we annotate as pre-LUCA or LUCA are
 144 contained within genes annotated by Moody et al. (21) as present in LUCA with more than 50%
 145 confidence, when present in their dataset (Figure 3c). This level of agreement far exceeds earlier
 146 works (22).



147

148 **Figure 3. Pfams (a) and clans (b) classified as ancient are well validated by the whole gene**
 149 **annotations of Moody et al. (21) (c).** a) Ancient post-LUCA Pfam classifications include 285
 150 LACA candidates and 2770 LBCA candidates (more analysis would be required to rule out
 151 extensive HGT within archaea or bacteria). Modern Pfams are distributed among the prokaryotic
 152 supergroups as follows: 9 CPR, 210 FCB, 942 Proteobacteria, 51 PVC, 1111 Terrabacteria, 2
 153 Asgard, 49 TACK, and 177 Euryarchaeota. In addition to supergroup-specific modern Pfams, we
 154 classified another 1097 Pfams, present in exactly two bacterial supergroups, as modern post-
 155 LBCA. We deemed 15 Pfams unclassifiable due to high inferred HGT rates, 397 due to
 156 uncertainty in rooting, and 198 due to ancient rooting combined with absence from too many
 157 supergroups (see Methods). b) Pre-LUCA clans contain at least two LUCA-classified Pfams or
 158 one pre-LUCA Pfam, whereas LUCA clans contain exactly one LUCA Pfam. Ancient post-LUCA
 159 clans contain no LUCA, pre-LUCA, or unclassified Pfams; they include an ancient post-LUCA
 160 Pfam or at least two modern Pfams covering at least two supergroups from only one of either
 161 bacteria or archaea. Modern clans include Pfams whose root is assigned at the origin of one

162 supergroup. Finally, unclassifiable clans did not meet any of our clan classification criteria, e.g.,
163 because they included both post-LUCA and unclassifiable Pfams. c) 98% of our pre-LUCA Pfams
164 and 87% of our LUCA Pfams are present in genes annotated by as present in LUCA with more
165 than 50% confidence, when present in their dataset. We mapped all Clusters of Orthologous
166 Genes (COGs) (95) in the Moody et al. (21) supplementary dataset (STable_1.csv) to UniProt IDs
167 (96) using the EggNOG 5.0 database (97). We then identified their associated Pfams using the
168 'Pfam-A.regions.uniprot.tsv' file downloaded from the Pfam FTP site ([https://pfam-](https://pfam-docs.readthedocs.io/en/latest/ftp-site.html#current-release)
169 [docs.readthedocs.io/en/latest/ftp-site.html#current-release](https://pfam-docs.readthedocs.io/en/latest/ftp-site.html#current-release)) (24) on May 28th, 2024. Our protein to
170 Pfam ID mappings are available in 'Protein2Domain_mappings' at
171 <https://doi.org/10.6084/m9.figshare.27191274.v1>.

172 In agreement with the Moody et al. (21) classification of LUCA metabolism, almost all Pfams
173 associated with enzymes in hydrogen metabolism, assimilatory nitrate and sulfate reduction
174 pathways, and the Wood-Ljungdahl pathway date back to LUCA (Supplementary Table 3). Our
175 results also support a, post-LUCA, bacterial origin of nitrogen fixation (21, 25) (Supplementary
176 Table 3). We assign to LUCA the complete set of amino acid-tRNA synthetase-associated anti-
177 codon binding domains found in modern prokaryotes. Here, focusing on complete genes would
178 have been problematic, because accessory amino acid-tRNA synthetase-associated domains
179 (e.g. PF04073 and PF13603, which deacylate misacylated tRNA) were sometimes added later.

180 We also checked the antiquity of the cofactor/cosubstrate S-adenosylmethionine (SAM) (26), both
181 with respect to SAM biosynthesis and SAM usage. In agreement with past work attributing the
182 SAM biosynthesis enzyme methionine adenosyltransferase to LUCA (27, 28), we assign its single
183 Pfam (PF01941) to LUCA (the corresponding COG1812 is not analyzed by Moody et al. (21)). In
184 agreement with past work attributing SAM-dependent methyltransferases to LUCA (29), Moody et
185 al. (21) assign the RsmB/RsmF family (COG0144), which methylates 16S rRNA, more than 75%
186 confidence of being present in LUCA, and we also classify its SAM-binding Rossmann fold Pfam
187 (PF01189) as LUCA. In agreement with (30, 31), Moody et al. (21) assign the SAM-binding tRNA
188 methylthiolase (COG0621) to LUCA with more than 75% confidence, and we confirm the pre-
189 LUCA status of its associated Radical SAM, TIM-barrel-related Pfam (PF04055). In agreement
190 with attribution of polyamines to LUCA (32) we assign to LUCA the one Pfam (PF02675) of S-
191 adenosylmethionine decarboxylase, which acts on SAM in the first step of polyamine synthesis;
192 the antiquity of corresponding COG1586 is not further confirmed by Moody et al. (21).

193

194 **Hydrophobic amino acids are more interspersed within ancient proteins**

195 Interspersion of hydrophobic amino acids away from one another along the primary sequence is
196 believed to mitigate risks from protein misfolding, while still enabling correct folding (33-35). Older
197 sequences have previously been found to have greater interspersion among their hydrophobic
198 residues, indicating more sophisticated protein folding (14, 36), likely due to survivorship bias
199 (37). Our Pfam age classifications confirm the antiquity of this trend, previously observed only for
200 animal sequences. LUCA Pfams show even more hydrophobic interspersion than the still-ancient
201 'post-LUCA' Pfams that include LACA candidates and LBCA candidates (Supplementary Figure
202 1; Wilcoxon rank sum test; $p = 0.02$). Post-LUCA Pfams in turn have more hydrophobic
203 interspersion than 'modern' Pfams that are specific to particular prokaryotic supergroups
204 (Wilcoxon rank sum test; $p = 0.02$).

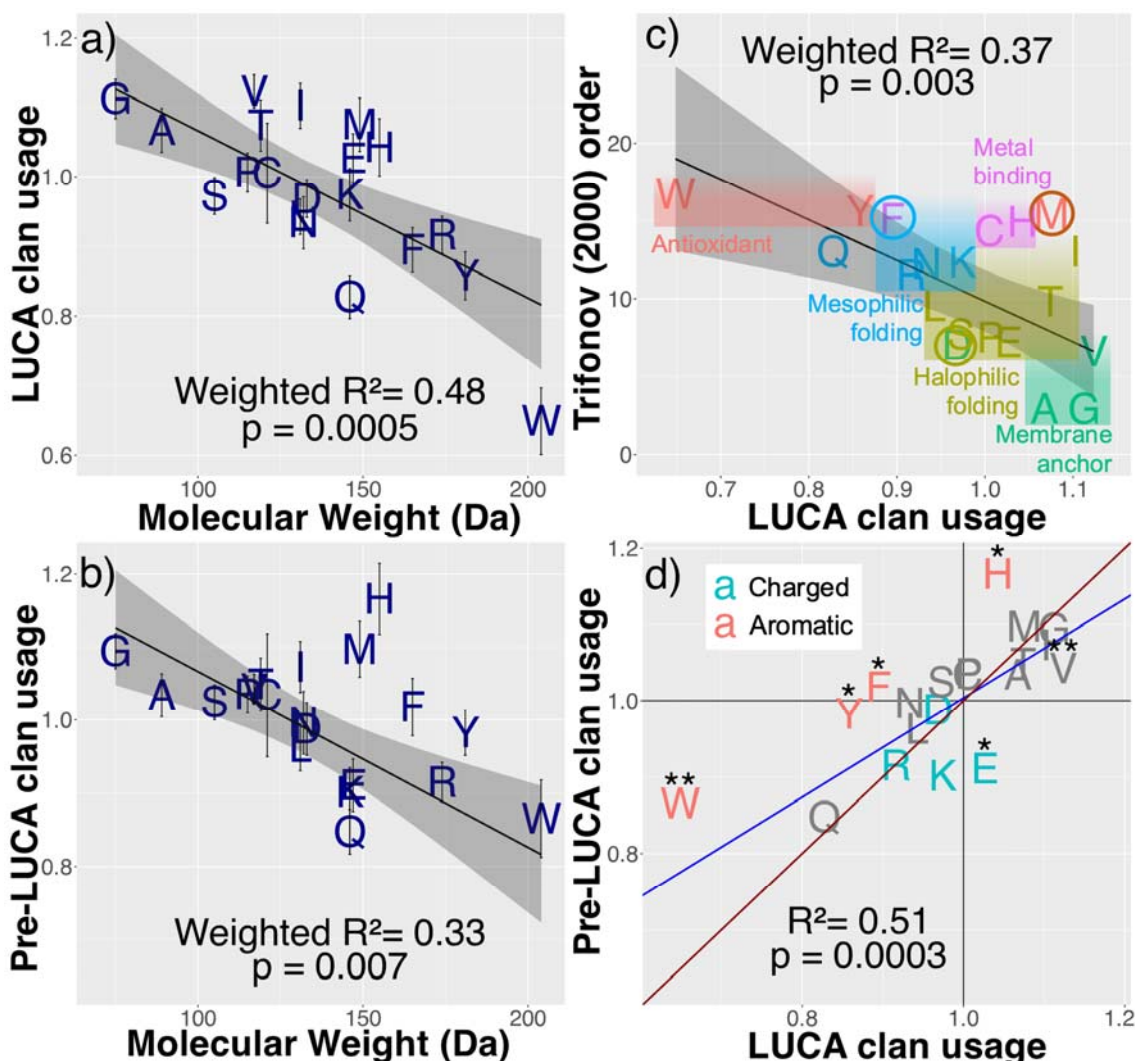
205

206 **LUCA's protein sequences were depleted in larger amino acids**

207 Clans present in LUCA were born before the divergence of Archaea and Bacteria, some
208 potentially prior to the completion of the genetic code. If newly recruited amino acids were added
209 slowly, the contemporary descendants of LUCA clans will show signs of ancestral depletion in
210 amino acids that were added late to the genetic code. We first focus on clans present in one copy

211 in LUCA (denoted “LUCA clans”), excluding those that had already duplicated and diverged into
 212 multiple surviving lineages (denoted “pre-LUCA clans”). We score ancestral amino acid
 213 enrichment and depletion as relative to still-ancient post-LUCA clans, which represent amino acid
 214 usage from the standard genetic code of all 20 amino acids, plus any ascertainment biases. This
 215 ratio, reflecting ancient amino acid usage, is not confounded with the effects of temperature, pH,
 216 oxygen tolerance, salinity, GC content, or transmembrane status on amino acid frequencies
 217 (Supplementary Figures 2a-f). Indeed, LUCA usage is similar in the very different biophysical
 218 context of a transmembrane site (Supplementary Figure 3).

219 Smaller amino acids are enriched in LUCA (Figure 4a; weighted $R^2 = 0.48$, $p = 0.0005$). Results
 220 are similar using a restricted set of Pfams validated by Moody et al. (21) (weighted $R^2 = 0.44$, $p =$
 221 0.001). As a negative control for methodological artifacts, the ancestral amino acid usage of post-
 222 LUCA clans relative to modern clans is not correlated with molecular weight ($p = 0.9$).



223

224 **Figure 4. LUCA is enriched for smaller amino acids, with subtle differences between single**
 225 **copy LUCA vs. multi-copy pre-LUCA sequences.** Ancestrally reconstructed amino acid
 226 frequencies in LUCA and pre-LUCA clans are shown relative to those in ancient post-LUCA
 227 clans. a) LUCA clans and b) pre-LUCA clans are enriched for amino acids of smaller
 228 weight. Weighted model 1 regression lines are shown in black with 95% confidence interval grey

229 shading. Error bars indicate standard errors. c) Character colors show the assignments of
230 Moosmann (5); colored circles indicate our re-assignments. We reclassify phenylalanine (F)
231 because it is enriched in proteins in mesophiles compared to their orthologs in thermophiles and
232 hyperthermophiles (98). We reclassify aspartic acid (D) because the surfaces of proteins within
233 halophilic bacteria are highly enriched in aspartic acid compared to in the surfaces of non-
234 halophilic mesophilic and thermophilic bacteria, in a manner that cannot be accounted for by the
235 dinucleotide composition of the halophilic genomes (99). The brown circle around methionine
236 highlights that while it might not be utilized against reactive oxygen species, it might once have
237 been against ancient reactive sulfur species. d) Model 2 Deming regression (accounting for
238 standard errors in both variables, implemented in deming() version 1.4-1 (100)) in blue shows that
239 pre-LUCA enrichments are not more extreme versions of LUCA enrichments, lying on the wrong
240 side of the $y=x$ red line. We include the imidazole-ring-containing H as aromatic. Asterisks (*)
241 indicate statistically different amino acid frequencies between pre-LUCA and LUCA (Welch two
242 sample t-test, $p < 0.05$ and $p < 0.01$).

243 **Revised Order of Amino Acid Recruitment**

244 Figure 4c visualizes how LUCA's amino acid enrichments compare to Trifonov's consensus order
245 (4). While they are correlated (weighted $R^2 = 0.37$, $p = 0.003$), this association disappears in a
246 weighted multiple regression with both molecular weight ($p = 0.03$) and Trifonov's (4) order ($p =$
247 0.9) as predictors (weighted $R^2 = 0.48$). This is also true using Trifonov's revised 2004 order
248 based on 60 metrics (38) (weighted $R^2 = 0.34$, $p = 0.006$ on its own; $p = 0.9$ when molecular
249 weight is also a predictor of LUCA usage). This suggests that some of Trifonov's 40-60 metrics
250 made his estimates of the order of recruitment worse rather than better. We use enrichment in
251 LUCA to re-classify VGIMTAHEPC as 'early' and depletion to classify KSDLNRFYQW as 'late'.
252 More precise estimation of the order of recruitment, with standard errors, is given in Table 1.

253 We place glutamine (Q or Gln) as the second last amino acid, much later than Trifonov (4)
254 inferred. Consistent with its late addition, Gln-tRNA synthetase (GlnRS) is either absent in
255 prokaryotes, or acquired via horizontal gene transfer from eukaryotes (39). Prokaryotes that lack
256 GlnRS perform tRNA-dependent amidation of Glu mischarged to Gln-tRNA by GluRS, forming
257 Gln-acylated Gln-tRNA via amidotransferase. The core catalytic domain (PF00587), shared
258 between the GlnRS and GluRS paralogs, is present in LUCA and can indiscriminately acylate
259 both Gln-tRNA and Glu-tRNAs with Glu (40).

260 **Metal-binding and sulfur-containing amino acids were added early to the genetic code**

261 Methionine (M), cysteine (C), and histidine (H) are all enriched in LUCA, despite previous
262 annotation as late additions to the genetic code (Figure 4c). C and H are the most frequently used
263 amino acids for binding iron, zinc, copper, and molybdenum, and H, aspartic acid (D) and
264 glutamic acid (E or Glu) for binding manganese and cobalt (Figure 2D of (41)). Binding can either
265 be to a metal ion, or to iron-sulfur (FeS) clusters, usually via C but sometimes via H or D (42).
266 Binding these transition metals is key to catalysis (43). Figure 4a is incompatible with C, H, D, or
267 E being late additions, and indeed H is more enriched than one would expect from its molecular
268 weight.

269 C and M are the only sulfur-containing amino acids in the contemporary genetic code.
270 Contemporary prokaryotes living in H_2S -rich environments use more C and M than matched
271 species (Supplementary Figure 4); LUCA's C and M enrichment might thus reflect an
272 environment rich in H_2S .

273 Moosmann (5) classified M, tryptophan (W), and tyrosine (Y) as antioxidants, because he
274 believed them to protect the overall protein structure from oxidative stress via sacrificial
275 oxidization. For instance, surface M residues can be reversibly oxidized to form methionine

276 sulfoxide (44). This might have driven isoleucine recoding to methionine in mitochondria (45, 46).
 277 However, proteins in aerobes are enriched in W and Y but not in M (47). Our results also
 278 separate early M from late Y and W (Figure 4). We speculate that methionine, abundant due to
 279 early life's use of SAM, might have protected against reactive sulfur species such as sulfide (S²⁻),
 280 which were present in early, H₂S-rich environments (48). Our results are then partially compatible
 281 with Granold et al.'s (49) view that Y and W (but not M) were added to complete the modern
 282 genetic code after reactive oxygen species became the main oxidizing threat.

Amino acid	LUCA usage	LUCA usage standard error	Pre-LUCA usage	Pre-LUCA usage standard error
V	1.12	0.0241	1.04	0.0205
G	1.11	0.0283	1.09	0.0241
I	1.1	0.0325	1.07	0.0351
M	1.08	0.0386	1.1	0.0383
A	1.07	0.0317	1.03	0.0297
T	1.07	0.0369	1.05	0.0362
H	1.04	0.0416	1.17	0.0486
E	1.03	0.0357	0.911	0.0357
C	1.01	0.0722	1.03	0.0844
P	1.01	0.0282	1.04	0.0255
K	0.974	0.038	0.901	0.0334
S	0.972	0.0265	1.02	0.0239
D	0.968	0.027	0.988	0.0363
L	0.942	0.0256	0.962	0.032
N	0.934	0.0374	0.996	0.0432
R	0.916	0.0265	0.915	0.0271
F	0.895	0.032	1.02	0.0394
Y	0.858	0.0341	0.982	0.0309
Q	0.827	0.031	0.847	0.0304
W	0.649	0.0476	0.865	0.0526

283 **Table 1.** LUCA and pre-LUCA clans' ancestral amino acid frequencies are divided by post-LUCA
 284 clan's ancestral amino acid frequencies to produce measures of relative usage. The standard
 285 errors of the amino acid usages were calculated using an approximation derived from a Taylor
 286 expansion of the ratio (90). For each of the 20 ancestral amino acid frequencies, the standard
 287 errors of the weighted means across all the clans within the LUCA and pre-LUCA phylostrata
 288 (weighted by the maximum number of ancestral sites across all Pfams in a given clan) were
 289 calculated using the weighted_se() function in the diagis R package (89)(See Methods for more
 290 detail).

291

292 Pre-LUCA clans hint at more ancient genetic codes

293 We expected pre-LUCA enrichments and depletions to be more extreme than for LUCA, but only
 294 H fits this prediction (Figure 4d), with significantly higher frequencies in pre-LUCA than in LUCA.

295 There is nevertheless a strong overall correlation between LUCA and pre-LUCA usages ($R^2 =$
296 0.51 , $p = 0.0003$). Pre-LUCA, like LUCA, is strongly depleted in Q, supporting the inference that
297 Q, not Y, was the 19th amino acid recruited into the standard genetic code. Pre-LUCA usage does
298 not correlate with Trifonov's consensus order (4) ($p = 0.2$), and correlates more weakly with
299 molecular weight (Figure 4b) (weighted $R^2 = 0.33$, $p = 0.007$).

300 H is one of six amino acids with significantly different frequencies in pre-LUCA vs. LUCA. All
301 three of the canonical, benzene-ring bearing, aromatic amino acids (W, Y, and phenylalanine (F)),
302 as well as the imidazole-ring containing H, are more common in pre-LUCA than in LUCA (Figure
303 4d, Welch 2-sample t-test; $p = 0.03$, 0.001 , 0.03 and 0.01 , respectively; 2.4% vs 2.1% H, 1.2% vs
304 0.9% W, 3.1% vs. 2.8% Y, and 4.1% vs. 3.7% F). Glutamic acid (E) and Valine (V) are less
305 common in pre-LUCA than in LUCA (Welch 2-sample t-test; $p = 0.01$ and 0.004 , respectively;
306 7.3% vs. 8.2% E, 7.5% vs. 8.1% V).

307 More W in pre-LUCA than LUCA is particularly surprising, because there is scientific consensus
308 that W was the last of the 20 canonical amino acids to be added to the genetic code. Therefore,
309 we manually inspected the pre-LUCA Pfam with the highest tryptophan frequency (3.1%):
310 PF00133, the core catalytic domain of the tRNA synthetases of leucine (L), isoleucine (I), and
311 valine (V). Each of these three synthetases has well-separated archaeal and bacterial branches,
312 confirming its pre-LUCA dating (Supplementary Figure 5). Highly conserved tryptophan sites
313 regulate the size of the amino acid binding pocket, allowing the synthetases to discriminate
314 among I, L, and V (50). There are also conserved I and V sites in the common ancestor of the I
315 and V tRNA synthetases, indicating that discrimination between the two happened prior to the
316 evolution of the synthetases currently responsible for the discrimination (51). This suggests that
317 an alternative, more ancient system predated the modern genetic code, and in particular predated
318 the evolution of super-specific, cognate aaRSs (51).

319

320 Discussion

321 The evolution of the current genetic code proceeded via stepwise incorporation of amino acids,
322 driven in part by changes in early life's environment and requirements. Contemporary proteins
323 retain information about which amino acids were part of the code at the time of their birth,
324 allowing us to infer the order of recruitment on the basis of enrichment or depletion in LUCA's
325 protein domains. Smaller amino acids were added to the code first, and when this is accounted
326 for, there is no further information in Trifonov's (4) widely used 'consensus' order based on 40
327 metrics, some of dubious relevance. The sulfur-containing amino acids C and M were
328 incorporated earlier than previously thought, likely because those metrics included experiments
329 conducted in the absence of sulfur. Q was added later than previously thought, in agreement with
330 evidence from glutamyl-tRNA synthetases. M and H were added to the code earlier than
331 expected from their molecular weights, and Q later. Even more ancient amino acid usage, in
332 sequences that had already duplicated and diverged pre-LUCA, shows significantly higher
333 frequencies of the aromatic amino acids W, Y, F, and H, and significantly lower frequencies of E
334 and V.

335 If LUCA lived in a H₂S-rich environment (48, 52), M residues could have protected proteins
336 against sulfur-mediated oxidative stress. M would furthermore have had high biotic availability as
337 the precursor (53) and product (54) of SAM, given our finding that LUCA made and used SAM.
338 The potentially sulfur-rich nature of early terrestrial life is context for astrobiology investigations of
339 sulfur-rich environments on Mars and Europa, with associated biosignatures key to life detection
340 (55).

341 An early role for H is compatible with a key role for metal binding in early life. It also resolves the
342 previous puzzle that the ancestral, conserved region of all Class I aaRSs contains a histidine-rich
343 HIGH motif (56, 57). The lack of abiotic availability was key to H's previous annotation as late, but
344 biotic availability of H in an RNA-dominant biotic context would have been sufficient. The

345 importance of abiotic availability (58, 59) to the origins of the genetic code remains unclear. We
346 note that ongoing research on plausible prebiotic syntheses in cyanosulfidic environments (60)
347 and alkaline hydrothermal vents (61) is reshaping our understanding of which amino acids were
348 accessible to early life. Amino acid abundances obtained from asteroid sample returns will also
349 soon contribute (62).

350 Our results offer an improved approximation of the order of recruitment of the twenty amino acids
351 into the genetic code under which contemporary protein-coding sequences were born. This order
352 need not match the importance or abundance with which amino acids were used by still earlier
353 life forms, nor during the prebiotic to biotic transition. Instead of using Trifonov's assignments (4)
354 to capture the order in which amino acids were recruited into our genetic code, we recommend
355 using the LUCA amino acid enrichment values plotted on the y-axis of Figure 4a, which can be
356 found together with their standard errors in Table 1.

357 More broadly, coding for different amino acids might have emerged at similar times but in
358 different biogeochemical environments. The temporal order of recruitment that we infer based on
359 LUCA sequences is not the temporal order for coding as a whole, but for the ancestor of the
360 modern translation machinery. Indeed, horizontal gene transfer of the tRNAs coupled with their
361 cognate aminoacyl tRNA synthetases might have brought the diverse components of the modern
362 translation machinery together (63). This further emphasizes that the time of origin of the
363 translation machinery's components need not match the time of their incorporation into the
364 surviving ancestral lineage.

365 The construction of the genetic code was tethered to the evolution of the ribosome (64). If the
366 ribosome's exit tunnel, whose formation and subsequent extension was key to ribosome evolution
367 (65), limited the size of the amino acids passing through, its progressive dilation could explain the
368 strong relationship between amino acid size and order of recruitment evidenced in LUCA clans. If
369 older, alternative codes were not similarly limited, this would explain why amino acid size is a
370 weaker predictor of pre-LUCA's amino acid usage compared to LUCA's amino acid usage.

371 To explain the different enrichments of pre-LUCA versus LUCA sequences, as well as the
372 surprising conservation of some sites prior to the emergence of the aaRSs that distinguish the
373 relevant amino acids, we propose that some pre-LUCA sequences are older than the current
374 genetic code, perhaps even tracing back to a peptide world at the dawn of precellular life (7).
375 Stepwise construction of the current code and competition among ancient codes could have
376 occurred simultaneously (66, 67). Ancient codes might also have used non-canonical amino
377 acids, such as norvaline and norleucine (68) which can be recognized by LeuRS (69, 70). Along
378 with having different genetic codes, we speculate that pre-LUCA and LUCA might have existed in
379 different geochemical settings. For instance, if pre-LUCA ancestors inhabited alkaline
380 hydrothermal vents, where abiotically produced aromatic amino acids have been found (61), this
381 would explain their enrichment in pre-LUCA relative to LUCA. We note that abiotic synthesis of
382 aromatic amino acids might be possible in the water-rock interface of Enceladus's subsurface
383 ocean, which is speculated to be analogous to terrestrial alkaline hydrothermal vents (71).

384 Perhaps the biggest mystery is how sequences such as the common ancestor of L/I/V-tRNA
385 synthetase, which were translated via alternative or incomplete genetic codes, ended up being re-
386 coded for translation by the direct ancestor of the canonical genetic code. Harmonization of
387 genetic codes facilitated innovation sharing via HGT, making it advantageous to use the most
388 common code, driving code convergence (72, 73). Only once a common code was established
389 did HGT drop to levels such that a species tree became apparent, i.e. the LUCA coalescence
390 point corresponds to convergence on a code (72). Our identification of pre-LUCA sequences
391 provides a rare source of data about early, alternative codes.

392

393 **Materials and Method**

394 **Pfam sequences**

395 We downloaded genomes of 3562 prokaryotic species from NCBI that were present in the Web of
396 Life (WoL): Reference phylogeny of microbes (74) in August 2022. We classified them into five
397 bacterial supergroups (FCB, PVC, CPR, Terrabacteria and Proteobacteria (75, 76)) and four
398 archaeal supergroups (TACK, DPANN, Asgard and Euryarchaeota (77, 78)). We included
399 incomplete genomes, to enhance coverage of underrepresented supergroups.

400 We assign ages not to whole proteins but to each of their protein domain constituents. We used
401 InterProScan (79) to identify instances of each Pfam domain (24). We excluded Pfams with fewer
402 than 50 instances across all downloaded genomes. We also excluded 9 Pfams marked “obsolete”
403 starting July 2023. Among the remaining 8282 Pfams, 2496 Pfams had more than 1000
404 instances. We downsampled these to balance representation across the two taxonomic domains
405 (archaea and bacteria). For instance, a Pfam with 2000 bacterial and 500 archaeal instances was
406 downsampled by retaining all 500 archaeal sequences plus a subset (randomly sampled without
407 replacement) of 500 bacterial sequences.

408 The Pfam database includes annotations of “clans” of Pfams that share a common ancestor
409 despite limited sequence similarity; for many analyses, we used clans rather than Pfams to
410 ensure independent datapoints. We treated Pfams that were not annotated as part of a clan as
411 single-entry clans, with clan ID equal to their Pfam ID.

412

413 **Pfam trees**

414 We aligned downsampled sequences for each Pfam using MAFFT v.7 (80), to infer a preliminary
415 tree with IQ-Tree (81), using a time non-reversible amino acid substitution matrix trained on the
416 Pfam database (NQ.PFAM) (82), and no rate heterogeneity among sites. Because most Pfams
417 are too short for reliable tree inference, we next reconciled preliminary Pfam trees with the WoL
418 prokaryotic species tree (74) using GeneRax (19). While there is no perfect species tree for
419 prokaryotes, reconciliation even with a roughly approximate tree can still provide benefits. We ran
420 GeneRax twice. The first run used the LG amino acid substitution model, a gamma distribution
421 with four discrete rate categories, and a Subtree Prune and Regraft (SPR) radius of 3. The
422 second run used the output of reconciled trees from the first run as input, and switched to an SPR
423 radius of 5, and the Q.PFAM amino acid substitution model (83), which was trained on the Pfam
424 dataset. We did not use NQ.PFAM, because time non-reversible models are only implemented in
425 IQ-Tree (82), and not in GeneRax. In both runs, we used the UndatedDTL probabilistic model to
426 compute the reconciliation likelihood. The second run of GeneRax reduced estimated transfer
427 rates by an additional 7% (Welch two sample t-test, $p = 10^{-12}$), indicating continued improvements
428 to the phylogenies.

429 We re-estimated the branch lengths of the reconciled Pfam trees in IQ-Tree using the NQ.PFAM
430 substitution model with no rate heterogeneity, then performed midpoint rooting using the phytools
431 R package (84) on these re-estimated branch lengths. As alternative rooting methods, we also
432 explored and rejected minimum variance (85), minimal ancestral deviation (86), and rootstraps
433 based on time non-reversible substitution models (87). The first two methods work best when
434 deviations from the molecular clock average out on longer timescales, which is not true for
435 phylogenies in which evolution e.g. at different temperatures causes sustained differences in
436 evolutionary rate. Indeed, minimum variance failed to resolve the prokaryotic supergroups as
437 separate clades, in visual inspection of PF00001, due to presumed genuine rate variation among
438 taxa. The latter produced very low confidence roots. In contrast, midpoint rooting largely
439 conformed to expectations for aaRSs once we implemented the procedure for outlier removal
440 described under “Classifying Pfam domains into ancient phylostrata” below.

441 We then implemented a new --enforce-gene-tree-root option in GeneRax, and ran GeneRax in
442 evaluation mode, with Q.PFAM+G as the substitution and rate heterogeneity models,
443 respectively. Evaluation mode re-estimates the reconciliation likelihood and the duplication,
444 transfer and loss (DTL) rates on a fixed tree, without initiating a tree search. Fifteen reconciled

445 Pfam trees had inferred transfer rates higher than 0.6, three times the seed transfer rate
446 implemented by GeneRax. We took this as a sign of poor tree quality, and annotated these 15
447 Pfams as of unclassifiable age.

448

449 **Filtering out HGT between archaea and bacteria**

450 Exclusion of horizontal gene transfer (HGT) between bacteria and archaea facilitates the
451 classification of a Pfam into LUCA (Figure 2a). To achieve this, we divided sequences into
452 “homogeneous groups”, meaning the largest monophyletic group in the Pfam tree for which the
453 corresponding species all belong to the same prokaryotic supergroup. Each homogeneous group
454 was considered as a candidate for exclusion, via its “focal node” separating it from its sister
455 group. To avoid over-pruning, we do not consider deep focal nodes that are 2 or fewer nodes
456 away from the root.

457 To be excluded, we first require the focal node to be ‘mixed’, meaning its descendants are found
458 within both Bacteria and Archaea. We next require the focal node to be labelled by GeneRax as
459 most likely a transfer (T), rather than a duplication (D) or speciation (S). Finally, to identify
460 homogeneous groups likely to be receivers rather than the donors of transferred sequences, we
461 require the sister lineage to contain no sequences present in the same supergroup as that
462 defining the homogeneous group in question. An example of filtering is shown in Figure 2b.

463 We ran the filtering process twice to address rare occasions of an intradomain HGT nested within
464 another intradomain HGT group. In the second filter, we apply the third criterion after pruning the
465 homogenous groups identified as HGT during the first filter.

466

467 **Classifying Pfam domains into ancient phylostrata**

468 We re-rooted the HGT-pruned Pfam trees using the midpoint.root function in the ‘phytools’ R
469 package (84), before classifying them into phylostrata (i.e. cohort of sequences of similar age).
470 Classification was based on the locations of the most recent common ancestors (MRCAs) of each
471 supergroup. For a LUCA Pfam, we require the root to separate the MRCAs of all bacterial
472 supergroups from the MRCAs of all archaeal supergroups (Figure 2a).

473 If there were no horizontal transfer, and the tree of a Pfam present in one copy in LUCA were
474 error-free, then the MRCAs for the nine supergroups would be two to four branches away from
475 the root. This is true even if our Pfam tree and/or species tree do not correctly capture the true
476 phylogenetic relationships among supergroups. However, we cannot ignore HGT; we did not filter
477 out the products of HGT between supergroups within Archaea or within Bacteria, only that of HGT
478 between Archaea and Bacteria. HGT from a more derived supergroup to a more basal
479 supergroup will move the inferred MRCA of the former further back in time. Given rampant HGT,
480 whether real or erroneously implied by Pfam tree error, we required Pfams to have their
481 supergroups’ MRCA two branches away from the root (Figure 2a).

482 Phylogenies with three or more basal bacterial supergroups and two or more basal archaeal
483 supergroups were classified as LUCA. In other words, we allow the absence of up to two
484 supergroups per taxonomic domain, as compatible with ancestral presence followed by
485 subsequent loss. Trees with three or more basal bacterial supergroups but fewer than two basal
486 archaeal supergroups, as well as trees with two or more basal archaeal supergroups but fewer
487 than three basal bacterial supergroups, were classified as ancient but post-LUCA. These are
488 candidate Pfams for the Last Bacterial Common Ancestor (LBCA) and the Last Archaeal
489 Common Ancestor (LACA) phylostrata, respectively, but the necessary HGT filtering for sufficient
490 confidence in this classification is beyond the scope of the current work. If only one basal
491 supergroup is present, then the Pfam is classified into the corresponding supergroup-specific
492 phylostratum, meaning it emerged relatively recently (modern post-LUCA). If two basal bacterial
493 supergroups (and no archaeal supergroups) were present, the Pfam was classified as post-LBCA

494 which was also considered modern post-LUCA (younger than LBCA but older than the
495 supergroup-specific phylostrata). The remaining Pfams were considered unclassifiable.

496 We also classify into a pre-LUCA phylostratum the subset of LUCA-classified Pfams for which
497 there is evidence that LUCA contained at least two copies that left distinct descendants. This is
498 motivated by the assumption that LUCA domains that were born earlier are more likely to have
499 duplicated and diverged prior to the archaeal-bacterial split (88). We require that both the nodes
500 that are only one branch from the root be classified as LUCA nodes. This means that each of
501 these nodes should, after HGT filtering: i) split a pure-bacterial lineage from a pure-archaeal
502 lineage, and ii) include as descendants at least three bacterial and two archaeal basal MRCAs no
503 more than two nodes downstream of the potential LUCA nodes (Figure 2c).

504 Assignment of a Pfam to a phylostratum is sensitive to the root's position. Midpoint rooting is
505 based on the longest distance between two extant sequences. A single inaccurately placed
506 sequence can yield an abnormally long terminal branch, upon which the root is then based. This
507 phenomenon was readily apparent upon manual inspection of rooted Pfam trees. To ensure the
508 robustness of age classifications to the occasional misplaced sequence, we removed the Pfam
509 instance with the longest root-to-tip branch length in each HGT-filtered tree as potentially faulty,
510 re-calculated the midpoint root, and then re-classified each Pfam. We repeated this for ten
511 iterations, then retained only those Pfams that were classified into the same phylostratum at least
512 7 out of 10 times. Our HGT filtering algorithm does not act on nodes near the root, making it
513 robust to small differences in root position; we therefore did not repeat the HGT-filtering during
514 these iterations.

515 We classified clans that contained at least two LUCA Pfams as pre-LUCA clans. Clans that
516 contained both ancient archaeal and ancient bacterial post-LUCA Pfams (i.e. candidate LACA
517 and LBCA Pfams) were classified as LUCA. Clans that contained at least two different archaeal
518 but no bacterial supergroup-specific Pfams, or three different bacterial supergroup-specific Pfams
519 but no archaeal supergroup-specific Pfams, were classified as ancient post-LUCA clans. Clans
520 that meet neither of these criteria, and that contain at least one unclassified Pfam, were
521 considered unclassifiable due to the possibility that the unclassified Pfam might be older than the
522 classified Pfams present in the clan. All other clans were assigned the age of their oldest Pfam.

523 For a more stringent analysis of amino acid usage, we restrict our Pfam dataset to those present
524 in proteins annotated by Moody et al. (21) as >75% likely to be in LUCA. We then re-classified
525 clan ages. Data on the likelihood of Pfams being present in LUCA, as annotated by Moody et al.
526 (21), can be found in 'MoodyPfams_probabilities.csv' on GitHub.

527

528 **Ancestral amino acid usages**

529 Ancestral sequence reconstruction (ASR) can introduce a variety of biases. ASR methods do not
530 resolve alignment gaps well, to infer indel evolution, instead inferring ancestral sequences far
531 longer than any contemporary descendant. To avoid bias among amino acids regarding which
532 contemporary sequences appear in the ancestral sequence more often than they should, we
533 retain only sites where more than 50% of the sequences contain an amino acid (i.e. no indel).
534 This ensures that no amino acid can be double counted.

535 For Pfams classified as pre-LUCA or LUCA, we require that a given site contain an amino acid
536 and not a gap in at least five bacterial sequences and five archaeal sequences. This additional
537 filter helps ensure that the ancestrally reconstructed sites were not inserted post-LUCA (even
538 when the Pfam itself dates back to LUCA). It also reduces the impact of any Pfams misclassified
539 as ancient on the inferred ancient amino acid usage.

540 Following these filters, we ran the remaining sites in each Pfam alignment (prior to HGT filtering)
541 through IQ-Tree with the -asr option, the NQ.PFAM substitution model, and R10 rate
542 heterogeneity. We then excluded low confidence sites from subsequent analyses, based on the

543 most likely amino acid having an ancestral probability estimate <0.4. Combined with the other two
544 filters described above, the concatenated sequence length for all four phylostrata (pre-LUCA,
545 LUCA, post-LUCA, and modern) fell by ~11%, presumably preferentially excluding rapidly
546 evolving sites to a similar degree in all four cases, such that amino acid exclusion biases cancel
547 out when ratios are taken.

548 We then summed over the amino acid probability distributions at each site at the deepest node,
549 and divided by the number of sites, to obtain per-Pfam estimated ancestral amino acid
550 frequencies. For each clan, we took the ancestral amino acid frequencies across Pfams,
551 weighted by the number of ancestral sites in the Pfams. For each phylostratum, we averaged
552 across clans, weighted by the maximum number of ancestral sites across all Pfams in a given
553 clan. We calculated a standard error associated with each phylostratum mean using the
554 `weighted_se()` function in the `diagis` R package (89).

555 We divided ancestral amino acid frequencies for the LUCA and pre-LUCA phylostrata by post-
556 LUCA ancestral amino acid frequencies to produce measures of relative usage. Standard errors
557 of each of these ratios L/P were calculated using an approximation derived from a Taylor

558 expansion of the ratio: $\sqrt{\frac{\sigma_L^2}{P^2} + \frac{L^2 \sigma_P^2}{P^4}}$ (90). These were used in weighted linear model 1
559 regressions, using the `lm()` function with the 'weights' argument in the 'stats' package in base R
560 (91). Uncertainty in the ancestral states arising over 4 billion years of evolution is expected to
561 bring values of L/P closer to one, without entirely erasing the signal. As a negative control for
562 bias, we calculate the relative amino acid usage of post-LUCA clans by dividing the ancestral
563 amino acid frequencies for post-LUCA clans by the ancestral amino acid frequencies for modern
564 clans.

565 Standard errors in Trifonov's (4) average rank reflect but underestimate uncertainty; we therefore
566 treat Trifonov's (4) rankings as the dependent variable and use its weights rather than errors on
567 L/P to weight the regression model in Figure 4c. Standard errors are not available for alternative
568 results based on Trifonov's 2004 order (38).

569

570 **Hydrophobic interspersion**

571 The degree to which hydrophobic are clustered vs. interspersed along the primary sequence was
572 calculated as a normalized index of dispersion for each Pfam instance (35). This metric uses the
573 ratio of the variance to the mean in the number of the most hydrophobic amino acids (leucine,
574 isoleucine, valine, phenylalanine, methionine, and tryptophan) within consecutive blocks of six
575 amino acids. The values of this index of dispersion were then normalized, to make them
576 comparable across Pfams with different lengths and hydrophobicities. In cases where the Pfam
577 length was not a multiple of 6, the average across all possible 6-amino acid frames was
578 computed, trimming the ends as needed. For additional details, see Foy et al. (36) or James et al.
579 (14). For each Pfam, we then took the average across all its instances (prior to downsampling
580 species).

581

582 **Transmembrane annotation**

583 We identified transmembrane sites within each Pfam using DeepTMHMM (92) on a consensus
584 sequence generated from the original multiple sequence alignments (prior to HGT filtering) using
585 the majority-rule `seq_consensus()` function in the R package 'bioseq' (93).

586

587

588 **Data and Code Availability**

589 Data files and R scripts used to generate the results and figures are available at
590 [sawsanwehbi/Pfam-age-classification GitHub repository](#). Pfam sequences, alignments, trees and
591 mappings to protein IDs are available at [https://figshare.com/projects/Pfam-age-classification-](https://figshare.com/projects/Pfam-age-classification-data/201630)
592 [data/201630](https://figshare.com/projects/Pfam-age-classification-data/201630).

593

594 **Acknowledgments**

595 We thank NASA [80NSSC24K0384] and the John Templeton Foundation [62220] for funding SW,
596 NM, and JM, Chan-Zuckerberg Initiative [EOSS4-0000000312] for funding BQM, the DFG [STA
597 860/6-2] for funding BM, the NSF [DEB-2333243] for funding AW, NM, and JM, the Arizona
598 NASA Space Grant Consortium [80NSSC20M0041] for funding NM, and the NIH [T32GM132008]
599 for funding AW. We thank Mike Barker, Alan Moses, and Elisa Tomat for helpful discussions, and
600 Cat Wolner for comments on the manuscript.

601 References

- 602 1. F. H. C. Crick, The origin of the genetic code. *Journal of Molecular Biology* **38**, 367-379
603 (1968).
- 604 2. J. T.-F. Wong, A Co-Evolution Theory of the Genetic Code. *Proceedings of the National
605 Academy of Sciences* **72**, 1909-1912 (1975).
- 606 3. S. L. Miller, A Production of Amino Acids Under Possible Primitive Earth Conditions.
607 *Science* **117**, 528-529 (1953).
- 608 4. E. N. Trifonov, Consensus temporal order of amino acids and evolution of the triplet
609 code. *Gene* **261**, 139-151 (2000).
- 610 5. B. Moosmann, Redox Biochemistry of the Genetic Code. *Trends in Biochemical Sciences*
611 **46**, 83-86 (2021).
- 612 6. W. Nitschke, S. E. McGlynn, E. J. Milner-White, M. J. Russell, On the antiquity of
613 metalloenzymes and their substrates in bioenergetics. *Biochimica et Biophysica Acta
614 (BBA) - Bioenergetics* **1827**, 871-881 (2013).
- 615 7. S. D. Fried, K. Fujishima, M. Makarov, I. Cherepashuk, K. Hlouchova, Peptides before
616 and during the nucleotide world: an origins story emphasizing cooperation between
617 proteins and nucleic acids. *Journal of The Royal Society Interface* **19** (2022).
- 618 8. E. T. Parker *et al.*, Primordial synthesis of amines and amino acids in a 1958 Miller H₂S-
619 rich spark discharge experiment. *Proceedings of the National Academy of Sciences* **108**,
620 5526-5531 (2011).
- 621 9. C. S. Foden *et al.*, Prebiotic synthesis of cysteine peptides that catalyze peptide ligation
622 in neutral water. *Science* **370**, 865-869 (2020).
- 623 10. C. Shen, L. Yang, S. L. Miller, J. Oró, Prebiotic synthesis of histidine. *Journal of
624 Molecular Evolution* **31**, 167-174 (1990).
- 625 11. A. Vázquez-Salazar, A. Becerra, A. Lazcano, Evolutionary convergence in the
626 biosyntheses of the imidazole moieties of histidine and purines. *PLOS ONE* **13**,
627 e0196349 (2018).
- 628 12. A. D. Goldman, B. Kacar, Cofactors are Remnants of Life's Origin and Early Evolution.
629 *Journal of Molecular Evolution* **89**, 127-133 (2021).
- 630 13. A. J. M. Ribeiro, J. D. Tyzack, N. Borkakoti, G. L. Holliday, J. M. Thornton, A global
631 analysis of function and conservation of catalytic residues in enzymes. *Journal of
632 Biological Chemistry* **295**, 314-324 (2020).
- 633 14. J. E. James *et al.*, Universal and taxon-specific trends in protein sequences as a function
634 of age. *Elife* **10** (2021).
- 635 15. D. J. Brooks, J. R. Fresco, A. M. Lesk, M. Singh, Evolution of Amino Acid Frequencies in
636 Proteins Over Deep Time: Inferred Order of Introduction of Amino Acids into the Genetic
637 Code. *Molecular Biology and Evolution* **19**, 1645-1655 (2002).
- 638 16. G. P. Fournier, J. P. Gogarten, Signature of a Primitive Genetic Code in Ancient Protein
639 Lineages. *Journal of Molecular Evolution* **65**, 425-436 (2007).
- 640 17. P. K. Keese, A. Gibbs, Origins of genes: "big bang" or continuous creation? *Proceedings
641 of the National Academy of Sciences* **89**, 9489-9493 (1992).
- 642 18. S. B. Van Oss, A.-R. Carvunis, De novo gene birth. *PLoS Genet.* **15**, e1008160 (2019).
- 643 19. B. Morel, A. M. Kozlov, A. Stamatakis, G. J. Szöllösi, GeneRax: A tool for species tree-
644 aware maximum likelihood based gene family tree inference under gene duplication,
645 transfer, and loss. (2019).
- 646 20. M. C. Weiss *et al.*, The physiology and habitat of the last universal common ancestor. *Nat
647 Microbiol* **1**, 16116 (2016).
- 648 21. E. R. R. Moody *et al.*, The nature of the last universal common ancestor and its impact on
649 the early Earth system. *Nature Ecology & Evolution* 10.1038/s41559-024-02461-1 (2024).
- 650 22. A. J. Crapitto, A. Campbell, A. Harris, A. D. Goldman, A consensus view of the proteome
651 of the last universal common ancestor. *Ecology and Evolution* **12** (2022).
- 652 23. Y. Wang, H. Zhang, H. Zhong, Z. Xue, Protein domain identification methods and online
653 resources. *Comput Struct Biotechnol J* **19**, 1145-1153 (2021).

- 654 24. J. Mistry *et al.*, Pfam: The protein families database in 2021. *Nucleic Acids Research* **49**,
655 D412-D419 (2021).
- 656 25. H. W. Pi *et al.*, Origin and Evolution of Nitrogen Fixation in Prokaryotes. *Mol Biol Evol* **39**
657 (2022).
- 658 26. P. Laurino, D. S. Tawfik, Spontaneous Emergence of S-Adenosylmethionine and the
659 Evolution of Methylation. *Angewandte Chemie International Edition* **56**, 343-345 (2017).
- 660 27. B. P. S. Chouhan, M. Gade, D. Martinez, S. Toledo-Patino, P. Laurino, Implications of
661 divergence of methionine adenosyltransferase in archaea. *FEBS Open Bio* **12**, 130-145
662 (2022).
- 663 28. C. Minici *et al.*, Structures of catalytic cycle intermediates of the *Pyrococcus furiosus*
664 methionine adenosyltransferase demonstrate negative cooperativity in the archaeal
665 orthologues. *Journal of Structural Biology* **210**, 107462 (2020).
- 666 29. P. Laurino *et al.*, An Ancient Fingerprint Indicates the Common Ancestry of Rossmann-
667 Fold Enzymes Utilizing Different Ribose-Based Cofactors. *PLoS Biology* **14**, e1002396
668 (2016).
- 669 30. A. D. Goldman, J. T. Beatty, L. F. Landweber, The TIM Barrel Architecture Facilitated the
670 Early Evolution of Protein-Mediated Metabolism. *J Mol Evol* **82**, 17-26 (2016).
- 671 31. P. Z. Kozbial, A. R. Mushegian, Natural history of S-adenosylmethionine-binding proteins.
672 *BMC Struct Biol* **5**, 19 (2005).
- 673 32. A. J. Michael, Polyamine function in archaea and bacteria. *Journal of Biological*
674 *Chemistry* **293**, 18693-18701 (2018).
- 675 33. R. Schwartz, S. Istrail, J. King, Frequencies of amino acid strings in globular protein
676 sequences indicate suppression of blocks of consecutive hydrophobic residues. *Protein*
677 *Science* **10**, 1023-1031 (2001).
- 678 34. E. Monsellier *et al.*, The Distribution of Residues in a Polypeptide Sequence Is a
679 Determinant of Aggregation Optimized by Evolution. *Biophysical Journal* **93**, 4382-4391
680 (2007).
- 681 35. A. Irbäck, C. Peterson, F. Potthast, Evidence for nonrandom hydrophobicity structures in
682 protein chains. *Proceedings of the National Academy of Sciences* **93**, 9533-9538 (1996).
- 683 36. S. G. Foy, B. A. Wilson, J. Bertram, M. H. J. Cordes, J. Masel, A Shift in Aggregation
684 Avoidance Strategy Marks a Long-Term Direction to Protein Evolution. *Genetics* **211**,
685 1345-1355 (2019).
- 686 37. J. E. James, P. G. Nelson, J. Masel, Differential Retention of Pfam Domains Contributes
687 to Long-term Evolutionary Trends. *Mol Biol Evol* **40** (2023).
- 688 38. E. N. Trifonov, The Triplet Code From First Principles. *Journal of Biomolecular Structure*
689 *and Dynamics* **22**, 1-11 (2004).
- 690 39. V. Lamour *et al.*, Evolution of the Glx-tRNA synthetase family: the glutamyl enzyme as
691 a case of horizontal gene transfer. *Proceedings of the National Academy of Sciences* **91**,
692 8670-8674 (1994).
- 693 40. J. Lapointe, L. Duplain, M. Proulx, A single glutamyl-tRNA synthetase aminoacylates
694 tRNAGlu and tRNAGln in *Bacillus subtilis* and efficiently misacylates *Escherichia coli*
695 tRNAGln1 in vitro. *Journal of Bacteriology* **165**, 88-93 (1986).
- 696 41. J. Li *et al.*, The Metal-binding Protein Atlas (MbPA): An Integrated Database for Curating
697 Metalloproteins in All Aspects. *J Mol Biol* **435**, 168117 (2023).
- 698 42. C. Vallières *et al.*, Iron-sulfur protein odyssey: exploring their cluster functional versatility
699 and challenging identification. *Metallomics* **16** (2024).
- 700 43. C. Andreini, I. Bertini, G. Cavallaro, G. L. Holliday, J. M. Thornton, Metal ions in biological
701 catalysis: from enzyme databases to general principles. *JBIC Journal of Biological*
702 *Inorganic Chemistry* **13**, 1205-1218 (2008).
- 703 44. R. L. Levine, L. Mosoni, B. S. Berlett, E. R. Stadtman, Methionine residues as
704 endogenous antioxidants in \square proteins. *Proceedings of the National Academy of Sciences*
705 **93**, 15036-15040 (1996).

- 706 45. A. Bender, P. Hajjeva, B. Moosmann, Adaptive antioxidant methionine accumulation in
707 respiratory chain complexes explains the use of a deviant genetic code in mitochondria.
708 *Proceedings of the National Academy of Sciences* **105**, 16496-16501 (2008).
- 709 46. M. Schindeldecker, B. Moosmann, Protein-borne methionine residues as structural
710 antioxidants in mitochondria. *Amino Acids* **47**, 1421-1432 (2015).
- 711 47. S. Vieira-Silva, E. P. C. Rocha, An Assessment of the Impacts of Molecular Oxygen on
712 the Evolution of Proteomes. *Molecular Biology and Evolution* **25**, 1931-1942 (2008).
- 713 48. A. Neubeck, F. Freund, Sulfur Chemistry May Have Paved the Way for Evolution of
714 Antioxidants. *Astrobiology* **20**, 670-675 (2020).
- 715 49. M. Granold, P. Hajjeva, M. I. Toşa, F.-D. Irimie, B. Moosmann, Modern diversification of
716 the amino acid repertoire driven by oxygen. *Proceedings of the National Academy of
717 Sciences* **115**, 41-46 (2018).
- 718 50. S. Fukai *et al.*, Structural Basis for Double-Sieve Discrimination of L-Valine from L-
719 Isoleucine and L-Threonine by the Complex of tRNA^{Val} and Valyl-tRNA Synthetase. *Cell*
720 **103**, 793-803 (2000).
- 721 51. G. P. Fournier, C. P. Andam, E. J. Alm, J. P. Gogarten, Molecular Evolution of Aminoacyl
722 tRNA Synthetase Proteins in the Early History of Life. *Origins of Life and Evolution of
723 Biospheres* **41**, 621-632 (2011).
- 724 52. K. R. Olson, Hydrogen sulfide, reactive sulfur species and coping with reactive oxygen
725 species. *Free Radical Biology and Medicine* **140**, 74-83 (2019).
- 726 53. J. Komoto, T. Yamada, Y. Takata, G. D. Markham, F. Takusagawa, Crystal Structure of
727 the S-Adenosylmethionine Synthetase Ternary Complex: A Novel Catalytic Mechanism
728 of S-Adenosylmethionine Synthesis from ATP and Met. *Biochemistry* **43**, 1821-1831
729 (2004).
- 730 54. M. A. Grillo, S. Colombatto, S-adenosylmethionine and its products. *Amino Acids* **34**,
731 187-193 (2008).
- 732 55. A. Moreras-Marti, M. Fox-Powell, C. R. Cousins, M. C. Macey, A. L. Zerkle, Sulfur
733 isotopes as biosignatures for Mars and Europa exploration. *Journal of the Geological
734 Society* **179**, jgs2021-2134 (2022).
- 735 56. G. Eriani, M. Delarue, O. Poch, J. Gangloff, D. Moras, Partition of tRNA synthetases into
736 two classes based on mutually exclusive sets of sequence motifs. *Nature* **347**, 203-206
737 (1990).
- 738 57. J. Douglas, R. Bouckaert, C. W. Carter, R. Wills, Peter, Enzymic recognition of amino
739 acids drove the evolution of primordial genetic codes. *Nucleic Acids Research* **52**, 558-
740 571 (2024).
- 741 58. P. G. Higgs, R. E. Pudritz, A Thermodynamic Basis for Prebiotic Amino Acid Synthesis
742 and the Nature of the First Genetic Code. *Astrobiology* **9**, 483-490 (2009).
- 743 59. N. Kitadai, S. Maruyama, Origins of building blocks of life: A review. *Geoscience
744 Frontiers* **9**, 1117-1153 (2018).
- 745 60. J. Fairchild, S. Islam, J. Singh, D.-K. Bučar, M. W. Powner, Prebiotically plausible
746 chemoselective pantetheine synthesis in water. *Science* **383**, 911-918 (2024).
- 747 61. B. Menez *et al.*, Abiotic synthesis of amino acids in the recesses of the oceanic
748 lithosphere. *Nature* **564**, 59-63 (2018).
- 749 62. D. S. C. Lauretta, Harold C. Jr; Grossman, Jeffrey N. ; Polit, Anjani T. ; the OSIRIS-REx
750 Sample Analysis Team, OSIRIS-REx Sample Analysis Plan -- Revision 3.0. (2023).
- 751 63. T. Froese, J. I. Campos, K. Fujishima, D. Kiga, N. Virgo, Horizontal transfer of code
752 fragments between protocells can explain the origins of the genetic code without vertical
753 descent. *Scientific Reports* **8** (2018).
- 754 64. J. C. Bowman, A. S. Petrov, M. Frenkel-Pinter, P. I. Penev, L. D. Williams, Root of the
755 Tree: The Significance, Evolution, and Origins of the Ribosome. *Chemical Reviews* **120**,
756 4848-4878 (2020).
- 757 65. A. S. Petrov *et al.*, Evolution of the ribosome at atomic resolution. *Proceedings of the
758 National Academy of Sciences* **111**, 10251-10256 (2014).

- 759 66. D. W. Morgens, A. R. O. Cavalcanti, An Alternative Look at Code Evolution: Using Non-
760 canonical Codes to Evaluate Adaptive and Historic Models for the Origin of the Genetic
761 Code. *Journal of Molecular Evolution* **76**, 71-80 (2013).
- 762 67. E. V. Koonin, A. S. Novozhilov, Origin and evolution of the genetic code: The universal
763 enigma. *IUBMB Life* **61**, 99-111 (2009).
- 764 68. C. Alvarez-Carreño, A. Becerra, A. Lazcano, Norvaline and Norleucine May Have Been
765 More Abundant Protein Components during Early Stages of Cell Evolution. *Origins of Life*
766 *and Evolution of Biospheres* **43**, 363-375 (2013).
- 767 69. Y. Tang, D. A. Tirrell, Attenuation of the editing activity of the Escherichia coli leucyl-tRNA
768 synthetase allows incorporation of novel amino acids into proteins in vivo. *Biochemistry*
769 **41**, 10635-10645 (2002).
- 770 70. I. Apostol *et al.*, Incorporation of Norvaline at Leucine Positions in Recombinant Human
771 Hemoglobin Expressed in Escherichia coli. *Journal of Biological Chemistry* **272**, 28980-
772 28988 (1997).
- 773 71. C. R. Glein, J. A. Baross, J. H. Waite, The pH of Enceladus' ocean. *Geochimica et*
774 *Cosmochimica Acta* **162**, 202-219 (2015).
- 775 72. K. Vetsigian, C. Woese, N. Goldenfeld, Collective evolution and the genetic code.
776 *Proceedings of the National Academy of Sciences* **103**, 10696-10701 (2006).
- 777 73. G. Sella, D. H. Ardell, The Coevolution of Genes and Genetic Codes: Crick's Frozen
778 Accident Revisited. *Journal of Molecular Evolution* **63**, 297-313 (2006).
- 779 74. Q. Zhu *et al.*, Phylogenomics of 10,575 genomes reveals evolutionary proximity between
780 domains Bacteria and Archaea. *Nature Communications* **10** (2019).
- 781 75. C. Rinke *et al.*, Insights into the phylogeny and coding potential of microbial dark matter.
782 *Nature* **499**, 431-437 (2013).
- 783 76. C. T. Brown *et al.*, Unusual biology across a group comprising more than 15% of domain
784 Bacteria. *Nature* **523**, 208-211 (2015).
- 785 77. B. J. Baker *et al.*, Diversity, ecology and evolution of Archaea. *Nat Microbiol* **5**, 887-900
786 (2020).
- 787 78. W.-S. Shu, L.-N. Huang, Microbial diversity in extreme environments. *Nature Reviews*
788 *Microbiology* **20**, 219-235 (2021).
- 789 79. P. Jones *et al.*, InterProScan 5: genome-scale protein function classification.
790 *Bioinformatics* **30**, 1236-1240 (2014).
- 791 80. K. Katoh, D. M. Standley, MAFFT multiple sequence alignment software version 7:
792 improvements in performance and usability. *Mol Biol Evol* **30**, 772-780 (2013).
- 793 81. B. Q. Minh *et al.*, IQ-TREE 2: New Models and Efficient Methods for Phylogenetic
794 Inference in the Genomic Era. *Mol Biol Evol* **37**, 1530-1534 (2020).
- 795 82. C. C. Dang *et al.*, nQMaker: Estimating Time Nonreversible Amino Acid Substitution
796 Models. *Syst Biol* **71**, 1110-1123 (2022).
- 797 83. B. Q. Minh, C. C. Dang, L. S. Vinh, R. Lanfear, QMaker: Fast and Accurate Method to
798 Estimate Empirical Models of Protein Evolution. *Syst Biol* **70**, 1046-1060 (2021).
- 799 84. L. J. Revell, phytools: an R package for phylogenetic comparative biology (and other
800 things). *Methods in Ecology and Evolution* **3**, 217-223 (2012).
- 801 85. U. Mai, E. Sayyari, S. Mirarab, Minimum variance rooting of phylogenetic trees and
802 implications for species tree reconstruction. *PLOS ONE* **12**, e0182238 (2017).
- 803 86. F. D. K. Tria, G. Landan, T. Dagan, Phylogenetic rooting using minimal ancestor
804 deviation. *Nature Ecology & Evolution* **1** (2017).
- 805 87. S. Naser-Khdour, B. Quang Minh, R. Lanfear, Assessing Confidence in Root Placement
806 on Phylogenies: An Empirical Study Using Nonreversible Models for Mammals.
807 *Systematic Biology* **71**, 959-972 (2022).
- 808 88. G. P. Fournier, E. J. Alm, Ancestral Reconstruction of a Pre-LUCA Aminoacyl-tRNA
809 Synthetase Ancestor Supports the Late Addition of Trp to the Genetic Code. *J Mol Evol*
810 **80**, 171-185 (2015).

- 811 89. J. Helske, diagis: Diagnostic Plot and Multivariate Summary Statistics of Weighted
812 Samples from Importance Sampling. <https://doi.org/10.32614/CRAN.package.diagis>
813 (2023).
- 814 90. H. Seltman, Approximations for mean and variance of a ratio.
815 <https://www.stat.cmu.edu/~hseltman/files/ratio.pdf> (2012).
- 816 91. R Core Team, R: A Language and Environment for Statistical Computing. [https://www.R-](https://www.R-project.org/)
817 [project.org/](https://www.R-project.org/) (2024).
- 818 92. J. Hallgren *et al.*, DeepTMHMM predicts alpha and beta transmembrane proteins using
819 deep neural networks. *bioRxiv* 10.1101/2022.04.08.487609, 2022.2004.2008.487609
820 (2022).
- 821 93. F. Keck, Handling biological sequences in R with the bioseq package. *Methods in*
822 *Ecology and Evolution* **11**, 1728-1732 (2020).
- 823 94. W. F. Doolittle, The nature of the universal ancestor and the evolution of the proteome.
824 *Current Opinion in Structural Biology* **10**, 355-358 (2000).
- 825 95. R. L. Tatusov, M. Y. Galperin, D. A. Natale, E. V. Koonin, The COG database: a tool for
826 genome-scale analysis of protein functions and evolution. *Nucleic Acids Research* **28**,
827 33-36 (2000).
- 828 96. The UniProt Consortium, UniProt: the universal protein knowledgebase. *Nucleic Acids*
829 *Research* **45**, D158-D169 (2016).
- 830 97. J. Huerta-Cepas *et al.*, eggNOG 5.0: a hierarchical, functionally and phylogenetically
831 annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids*
832 *Research* **47**, D309-D314 (2018).
- 833 98. A. Szilágyi, P. Závodszy, Structural differences between mesophilic, moderately
834 thermophilic and extremely thermophilic protein subunits: results of a comprehensive
835 survey. *Structure* **8**, 493-504 (2000).
- 836 99. S. Fukuchi, K. Yoshimune, M. Wakayama, M. Moriguchi, K. Nishikawa, Unique amino
837 acid composition of proteins in halophilic bacteria. *J Mol Biol* **327**, 347-357 (2003).
- 838 100. T. Therneau, deming: Deming, Theil-Sen, Passing-Bablock and Total Least Squares
839 Regression. <https://doi.org/10.32614/CRAN.package.deming> (2024).
- 840