

MIT Open Access Articles

Paleoproterozoic sterol biosynthesis and the rise of oxygen

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

Citation: Gold, David A. et al. "Paleoproterozoic sterol biosynthesis and the rise of oxygen." Nature 543, 7645 (March 2017): 420–423. © 2017 Macmillan Publishers Limited, part of Springer Nature

As Published: <http://dx.doi.org/10.1038/nature21412>

Publisher: Springer Science and Business Media LLC

Persistent URL: <https://hdl.handle.net/1721.1/128450>

Version: Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

Terms of use: Creative Commons Attribution-Noncommercial-Share Alike



1 **Paleoproterozoic sterol biosynthesis and the rise of oxygen**

2

3 ^{1,2} David A. Gold (dagold@mit.edu)

4 ¹Abigail Caron (abigailc@mit.edu)

5 ¹Gregory P. Fournier (g4nier@mit.edu)

6 ¹Roger E. Summons* (rsummons@mit.edu)

7

8 1. Department of Earth, Atmospheric and Planetary Science, Massachusetts Institute
9 of Technology, Cambridge, MA, USA.

10 2. Present Address: Division of Biology and Biological Engineering, California
11 Institute of Technology, Pasadena, CA, USA.

12

13 * Corresponding author: rsummons@mit.edu

14

15

16 **Text**

17 Natural products preserved in the geologic record can function as “molecular fossils”,
18 providing insight into organisms and physiologies that existed in the deep past. One
19 important group of molecular fossils is the steroidal hydrocarbons (steranes), which are
20 the diagenetic remains of sterol lipids. Complex sterols with modified side chains are
21 unique to eukaryotes, although simpler sterols can also be synthesized by a few bacteria ¹.
22 Sterol biosynthesis is an oxygen-intensive process; thus, the presence of complex steranes
23 in ancient rocks not only signals the presence of eukaryotes, but also aerobic metabolic
24 processes ². In 1999, steranes were reported in 2.7 billion year old (Gyr-old) rocks from
25 the Pilbara Craton in Australia ³, suggesting a long delay between photosynthetic oxygen
26 production and its accumulation in the atmosphere (also known as the Great Oxidation
27 Event) 2.45–2.32 Gyr ago ⁴. However, the recent reappraisal and rejection of these
28 steranes as contaminants ⁵ pushes the oldest reported steranes forward to ~1.64 Gyr ⁶. In
29 the present study, we use a molecular clock approach to improve constraints on the
30 evolution of sterol biosynthesis. We infer that stem-eukaryotes shared functionally
31 modern sterol biosynthesis genes with bacteria *via* horizontal gene transfer. Comparing
32 multiple molecular clock analyses, we find that the maximum marginal probability for
33 the divergence time of bacterial and eukaryal sterol biosynthesis genes is ~2.31 Gyr ago,
34 concurrent with the most recent geochemical evidence for the Great Oxidation Event
35 (GOE) ⁷. Our results therefore suggest that simple sterol biosynthesis existed well before
36 the diversification of living eukaryotes, substantially predating the oldest detected sterane
37 biomarkers (~1.64 Gyr ⁶), and furthermore, that the evolutionary history of sterol

38 biosynthesis is tied to the first widespread availability of molecular oxygen in the ocean-
39 atmosphere system.

40

41 For this study we focused on the first two enzymes necessary for sterol biosynthesis:
42 squalene monooxygenase (SQMO; enzyme commission number 1.14.14.17) and
43 oxidosqualene cyclase (OSC; EC number 5.4.99.7). These proteins use molecular oxygen
44 to sequentially convert squalene into the protosterol precursors required for complex
45 eukaryotic sterols;⁸ they are also the only two enzymes in the canonical sterol
46 biosynthesis pathway (KEGG reference map00100) typically found in bacteria other than
47 squalene synthetase, which is broadly conserved across all domains of life. We wrote a
48 set of scripts (detailed in the Methods) to vet all potential SQMO and OSC proteins in the
49 NCBI protein database, map their distribution across taxa, and subsample the data based
50 on taxonomic completeness. To root these gene trees for phylogenetic analyses, we also
51 collected data from relevant protein paralogs: rooting OSC with squalene hopene cyclase
52 (SHC; EC number 5.4.99.17), and rooting SQMO with the ubiH/COQ6 family, which
53 includes the eukaryote-specific ubiquinone biosynthesis monooxygenase (COQ6), as well
54 as the bacterial ubiquinone biosynthesis hydroxylase family (UbiH/UbiF/VisC). This
55 protocol ensured a complete collection of prokaryotic SQMO and OSC proteins in NCBI,
56 as well as a well-sampled eukaryotic dataset.

57

58 Phylogenetic analyses of *sqmo* and *osc* suggest that the two genes have a shared
59 evolutionary history (Fig. 1a). The eukaryotic portions of the *sqmo* and *osc* gene trees
60 broadly recapitulate the most current species trees^{9,10} albeit with low support at the

61 deepest nodes (see Extended Data Figs. 1-6). This supports the hypothesis that both genes
62 were present in the last common ancestor of eukaryotes⁸. Some sequences from
63 amoebozoan, alveolate, and rhizarian taxa have unconventional placements within these
64 phylogenies; this is likely caused by long-branch attraction and other phylogenetic
65 artifacts, but we cannot rule out limited horizontal gene transfer between eukaryotes (see
66 Extended Data Figs. 1 and 2). We found no evidence for *sqmo* or *osc* genes in Archaea,
67 but did confirm one or both genes in 27 bacterial taxa, representing 6 phyla, 9 classes,
68 and 9 orders (Fig. 1b). In phylogenetic analyses, all bacterial sequences group in one of
69 two places: either basal to crown-group eukaryotes (“Bacterial Group I”), or nested
70 within the bikonts (“Bacterial Group II”). Similar trees were recently reported in a study
71 of bacterial *osc* genes¹, but the consistency between *osc* and *sqmo* topologies has not
72 previously been observed. The consistent grouping of bacterial taxa into either Group I or
73 II suggests that *sqmo* and *osc* have moved together through horizontal gene transfer (Fig.
74 1a). The linked inheritance of these two genes is also supported by synteny analysis of
75 bacterial genomes, with *sqmo* and *osc* rarely being separated by more than one gene (Fig.
76 1b). Such synteny has previously been observed in two species of bacteria¹¹, but this
77 study suggests the phenomenon is broadly conserved. The only taxon where *sqmo* and
78 *osc* demonstrate conflicting phylogenic placement is *Eudora adriatica*; because genomic
79 synteny is conserved in this bacterium, the inconsistency is unlikely to represent
80 independent gene transfers, and is most likely a long-branch artifact.

81

82 The distribution of *sqmo* and *osc* genes can be explained by at least two horizontal gene
83 transfer events between bacteria and eukaryotes. The nesting of Bacterial Group II within

84 the crown-group bikonts strongly suggests a horizontal gene transfer from eukaryotes
85 into bacteria. Conversely, Bacterial Group I roots outside of extant Eukarya, so the
86 directionality of this horizontal gene transfer event cannot be unambiguously determined
87 from our data by means of a polarizing outgroup. However, the shallow clade depth and
88 sparse phylogenetic distribution of bacterial genes—combined with the relatively long
89 branches leading to eukaryal representatives—are most consistent with these genes being
90 transferred one or more times from stem-eukaryotes to bacteria ¹¹. Additionally, because
91 Bacterial Group I includes genes from *Gemmata*—a genus that has been demonstrated to
92 use SQMO and OSC enzymes to produce protosterols ¹¹—we can infer that these
93 horizontal gene transfer events involved the sharing of functionally modern proteins,
94 permitting the reconstruction of the character state for these genes at the coalescent node.
95 In this way, horizontal gene transfer events can be used to map the presence of characters
96 onto stem lineages in the absence of paleontological evidence.

97

98 Insert Figure 1 hereabouts

99

100 To study the timing of the Bacterial Group I / stem-eukaryote split, we performed a series
101 of molecular clock analyses (Fig. 2 and Table 1). We tested SQMO and OSC
102 separately—with and without sister genes as outgroups—as well as the two proteins
103 concatenated into a single dataset. Our gene trees suggest that excavates are the earliest
104 branching eukaryote lineage, a hypothesis that is currently controversial ^{9,10}. We therefore
105 repeated all analyses using an alternative topology where excavates are sister to bikonts.
106 The timing of the Bacterial Group I / stem-eukaryote split has large uncertainties
107 associated with it—as expected given the limited data in single-gene analyses. The 95%

108 confidence intervals consistently fall between 1.75 and 3.05 Gyr, with one exception
109 coming from our SQMO-only datasets, which give significantly older dates than other
110 analyses (Table 1). But when SQMO is constrained by employing the UbiH/COQ6
111 family as an outgroup, it produces results congruent with the other analyses. The SQMO-
112 only datasets also produce an origin for crown-group eukaryotes that is significantly older
113 than estimates from previous multi-gene molecular clock analyses¹²⁻¹⁴, suggesting that
114 molecular clocks derived from SQMO-only data lead to a general overestimation of true
115 divergence times.

116

117

118 Insert Figure 2 hereabouts

119

120 The marginal probabilities associated with the Bacterial Group I / stem-eukaryote split
121 are greatest around the time period of the Great Oxidation Event (Fig. 3). The marginal
122 probability curves in the SQMO-only datasets are significantly older than the others, but
123 also have the lowest peak densities, again suggesting that these analyses should be
124 viewed with caution. Using our preferred analysis (a concatenated SQMO and OSC
125 sequences with an excavates-basal tree), we specifically recover a 94.5% probability of
126 this younger-bound age constraint on oxygen-dependent sterol biosynthesis predating the
127 Orosirian Period (>2.05 Ga; see Extended Table 1 for the distribution of marginal
128 probabilities for all analyses across geologic time). If we treat excavates as sister to the
129 bikonts (instead of sister to all eukaryotes, as our analysis supports), we recover slightly
130 younger dates, but the change does not fundamentally impact our interpretation of the
131 data. All analyses suggest a Paleoproterozoic-or-earlier existence for *sqmo* and *osc* genes.

132

133 Insert Figure 3 hereabouts

134

135 In conclusion, our molecular clock analyses suggest that protosterol biosynthesis likely
136 existed by the time oxygen is first detectable as a permanent presence in the Earth's
137 atmosphere^{4,7,15,16}. Our results similarly suggest that sterol biosynthesis substantially
138 predates the evolution of crown-group eukaryotes, and was likely an important
139 preadaptation to modern eukaryotic life¹⁷. Finally, our results are inconsistent with an
140 origin of sterol biosynthesis ~1.64 Gyr ago, as informed by the oldest sterane biomarkers
141 currently known⁶. Reconciling the molecular clocks with the Proterozoic biomarker
142 record requires considerable caution following the disproving of earlier reports of
143 Archaean steranes⁵, and the demonstration of widespread fossil fuel contamination in
144 laboratory aerosols¹⁸. Most reliance can be placed on studies conducted on low maturity,
145 organic-rich rock sequences and which have been replicated^{6,19} and/or supported using
146 multiple technologies²⁰. The progression of sterols identified in such studies—from
147 simple steroids at ~1.64 Gyr ago⁶, to atypical triterpenoid biomarker patterns in the
148 Neoproterozoic^{19,21}, to unambiguous and abundant eukaryotic sterols with modified side-
149 chains in the Phanerozoic—is consistent with the genetic data, but molecular clocks
150 suggest this progression must have occurred much earlier. Regarding Phanerozoic sterols
151 with modified side-chains, the conservation of the sterol biosynthesis pathway across
152 eukaryotes^{8,22} means that such sterols were being synthesized by the ancestral crown-
153 group eukaryote, which we date between ~1.30-2.17 Gyr ago, and multi-gene molecular
154 clocks estimate between 0.95 and 1.87 Gyr ago¹²⁻¹⁴. The presence of atypical sterane

155 patterns in 0.8-1.64 Gyr-old rocks is therefore consistent with crown-group eukaryotic
156 sources, although we cannot rule out non-eukaryotic interpretations for some of these
157 biomarkers^{19,21,23}. Regarding simple steroids, the most conservative interpretation of our
158 data (i.e. the youngest date falling within a 95% confidence interval in any of our
159 analyses) suggests that protosterols were being synthesized >1.75 Gyr ago. Using more
160 realistic estimates (based on the averaged means of all analyses, excluding the SQMO-
161 only datasets as outliers) we find that basic protosterols were likely being synthesized
162 >2.31 Gyr ago. This suggests a >670 Myr gap between our age estimates and the oldest
163 fossil steranes. This gap between the molecular and geochemical evidence for sterol
164 biosynthesis, which may reflect sampling bias or could have ecological²⁴ or taphonomic
165 explanations⁵, will only be resolved by further discovery (see Supplemental Text for a
166 more detailed discussion).

167

168 **References**

- 169 1 Wei, J. H., Yin, X. & Welander, P. V. Sterol Synthesis in Diverse Bacteria.
170 *Frontiers in Microbiology* **7**, 990, doi:10.3389/fmicb.2016.00990 (2016).
- 171 2 Summons, R. E., Bradley, A. S., Jahnke, L. L. & Waldbauer, J. R. Steroids,
172 triterpenoids and molecular oxygen. *Philosophical Transactions of the Royal*
173 *Society B-Biological Sciences* **361**, 951-968, doi:10.1098/rstb.2006.1837|ISSN
174 0962-8436 (2006).
- 175 3 Brocks, J. J., Logan, G. A., Buick, R. & Summons, R. E. Archean Molecular
176 Fossils and the Early Rise of Eukaryotes. *Science* **285**, 1033-1036,
177 doi:10.1126/science.1096806 (1999).
- 178 4 Bekker, A. *et al.* Dating the rise of atmospheric oxygen. *Nature* **427**, 117-120
179 (2004).
- 180 5 French, K. L. *et al.* Reappraisal of hydrocarbon biomarkers in Archean rocks.
181 *Proceedings of the National Academy of Sciences* **112**, 5915-5920 (2015).
- 182 6 Brocks, J. J. *et al.* Biomarker evidence for green and purple sulphur bacteria in a
183 stratified Palaeoproterozoic sea. *Nature* **437**, 866 (2005).
- 184 7 Luo, G. *et al.* Rapid oxidation of Earth's atmosphere 2.33 billion years ago.
185 *Science Advances* **2**, e1600134, doi:10.1126/sciadv.1600134 (2016).

- 186 8 Desmond, E. & Gribaldo, S. Phylogenomics of Sterol Synthesis: Insights into the
187 Origin, Evolution, and Diversity of a Key Eukaryotic Feature. *Genome Biology*
188 *and Evolution* **1**, 364-381, doi:10.1093/gbe/evp036 (2009).
- 189 9 He, D. *et al.* An Alternative Root for the Eukaryote Tree of Life. *Current Biology*
190 **24**, 465-470, doi:http://dx.doi.org/10.1016/j.cub.2014.01.036 (2014).
- 191 10 Derelle, R. *et al.* Bacterial proteins pinpoint a single eukaryotic root. *Proceedings*
192 *of the National Academy of Sciences of the United States of America* **112**, E693-
193 E699 (2015).
- 194 11 Pearson, A., Budin, M. & Brocks, J. J. Phylogenetic and biochemical evidence for
195 sterol synthesis in the bacterium *Gemmata obscuriglobus*. *Proceedings of the*
196 *National Academy of Sciences of the United States of America* **100**, 15352-15357
197 (2003).
- 198 12 Douzery, E. J. P., Snell, E. A., Baptiste, E., Delsuc, F. & Philippe, H. The timing
199 of eukaryotic evolution: Does a relaxed molecular clock reconcile proteins and
200 fossils? *Proceedings of the National Academy of Sciences of the United States of*
201 *America* **101**, 15386-15391 (2004).
- 202 13 Berney, C. & Pawlowski, J. A molecular time-scale for eukaryote evolution
203 recalibrated with the continuous microfossil record. *Proceedings of the Royal*
204 *Society of London B: Biological Sciences* **273**, 1867-1872 (2006).
- 205 14 Parfrey, L. W., Lahr, D. J. G., Knoll, A. H. & Katz, L. A. Estimating the timing of
206 early eukaryotic diversification with multigene molecular clocks. *Proceedings of*
207 *the National Academy of Sciences* **108**, 13624-13629 (2011).
- 208 15 Farquhar, J., Bao, H. & Thiemens, M. Atmospheric Influence of Earth's Earliest
209 Sulfur Cycle. *Science* **289**, 756-758, doi:10.1126/science.289.5480.756 (2000).
- 210 16 Lyons, T. W., Reinhard, C. T. & Planavsky, N. J. The rise of oxygen in Earth's
211 early ocean and atmosphere. *Nature* **506**, 307-315, doi:10.1038/nature13068
212 (2014).
- 213 17 Cavalier-Smith, T. The phagotrophic origin of eukaryotes and phylogenetic
214 classification of protozoa. *International Journal Of Systematic And Evolutionary*
215 *Microbiology* **52**, 297-354 (2002).
- 216 18 Illing, C. J., Hallmann, C., Miller, K. E., Summons, R. E. & Strauss, H. Airborne
217 hydrocarbon contamination from laboratory atmospheres. *Organic Geochemistry*
218 **76**, 26-38, doi:http://dx.doi.org/10.1016/j.orggeochem.2014.07.006 (2014).
- 219 19 Summons, R. E. *et al.* Distinctive hydrocarbon biomarkers from fossiliferous
220 sediment of the Late Proterozoic Walcott Member, Chuar Group, Grand Canyon,
221 Arizona. *Geochimica et Cosmochimica Acta* **52**, 2625 (1988).
- 222 20 Love, G. D., Stalvies, C., Grosjean, E., Meredith, W. & Snape, C. E. Analysis of
223 molecular biomarkers covalently bound within Neoproterozoic sedimentary
224 kerogen in *Paleontological Society Papers* Vol. 14 (eds P.H. Kelley & R.K.
225 Bambach) 67-83 (The Paleontological Society, 2008).
- 226 21 Brocks, J. J. *et al.* Early sponges and toxic protists: possible sources of cryostane,
227 an age diagnostic biomarker antedating Sturtian Snowball Earth. *Geobiology* **14**,
228 129-149, doi:10.1111/gbi.12165 (2016).
- 229 22 Gold, D. A. *et al.* Sterol and genomic analyses validate the sponge biomarker
230 hypothesis. *Proceedings of the National Academy of Sciences* **113**, 2684-2689
231 (2016).

232 23 Banta, A. B., Wei, J. H. & Welander, P. V. A distinct pathway for tetrahymanol
233 synthesis in bacteria. *Proceedings of the National Academy of Sciences* **112**,
234 13478-13483 (2015).

235 24 Anbar, A. D. & Knoll, A. H. Proterozoic ocean chemistry and evolution: A
236 bioinorganic bridge? *Science* **297**, 1137-1142 (2002).

237

238 **Legends to Figures**

239

240 Figure 1. Phylogeny and synteny of *sqmo* and *osc* genes. (A) Maximum likelihood trees
241 for *SQMO* (left) and *OSC* (right) based on Extended Data Figs. 3 and 4. Note that the two
242 trees provide consistent topologies, with bacteria clustering within crown-group
243 eukaryotes (“Bacterial Group II”) or outside (“Bacterial Group I”). Some eukaryotes have
244 been excluded from these trees; see Extended Data Figs. 1 and 2. (B) Distribution of
245 vetted bacterial *SQMO* and *OSC* genes, divided into the two phylogenetic groups. For
246 species containing both genes, the relevant placement of both genes on the genome is
247 provided to the left. Additional genes of interest have also been colored, and are
248 described in the key. Note how there is greater evidence of synteny in Bacterial Group II,
249 consistent with this clade representing a more recent horizontal gene transfer event.

250

251 Figure 2. Molecular clock for one of the datasets used in this study (*SQMO* + *OSC*
252 concatenated together, with excavates as the earliest-branching eukaryotes). Bars
253 represent 95% confidence intervals for nodes with >70% posterior probability. See Table
254 1 for the results of all analyses. Stars signify important nodes included in Table 1, circles
255 indicate fossil calibrations.

256

257 Figure 3. Marginal probability curves for the timing of the Bacterial Group I / stem-
258 eukaryote split (red star in Figure 2). Differences in clock estimation times are compared
259 for our preferred topology (A), where excavates are the earliest-branching eukaryotes,
260 and (B) where excavates are sister to the bikonts. Curves were generated by sampling the
261 MCMC analysis every 1,000 generations for 20 million generations, with a 25% burn-in.
262

263 **Acknowledgements**

264 We gratefully acknowledge funding from the Agouon Institute Geobiology Fellowship
265 to D.A.G. and the Simons Foundation Collaboration on the Origins of Life (SCOL) to
266 R.E.S. Additional support was provided by the NSF program ‘Frontiers of Earth System
267 Dynamics’ (EAR-1338810) to R.E.S.

268

269 **Author contributions**

270 R.E.S. and D.A.G. designed the experiment. D.A.G. and A.M.C. performed the data
271 analysis. All authors were involved in interpreting the data and drafting the manuscript.

272

273 **Competing financial interests.**

274 The authors declare no competing financial interests.

275

276 **Materials & Correspondence.**

277 Correspondence to: Roger Summons (rsummons@mit.edu)

278

279

280 **Table**

281 Table 1. Mean values for important nodes in molecular clock analyses (in Gyr), with 95%

282 confidence intervals in parentheses.

<u>Gene</u>	<u>Gene outgroup included?</u>	<u>Excavates basal eukaryotes?</u>	<u>Origin of Bacterial Group I</u>	<u>Origin of crown-eukaryote SQMO/OSC</u>	<u>Origin of Bacterial Group II</u>
SQMO	No	Yes	3.21(2.49-4.03)	2.22 (1.78-2.73)	1.54 (1.30-1.83)
SQMO	No	No	3.27 (2.51-4.18)	1.98 (1.58-2.34)	1.61 (1.31-1.92)
SQMO	Yes	Yes	2.47 (1.98-3.05)	1.80 (1.51-2.12)	1.35 (1.22-1.50)
SQMO	Yes	No	2.34(1.78-2.95)	1.60(1.40-1.92)	1.36(1.23-1.53)
OSC	No	Yes	2.23(1.78-2.77)	1.61(1.39-1.87)	1.34(1.15-1.57)
OSC	No	No	2.23(1.75-2.75)	1.51(1.30-1.81)	1.34(1.10-1.64)
OSC	Yes	Yes	2.34(1.93-2.82)	1.56(1.37-1.78)	1.32(1.14-1.51)
OSC	Yes	No	2.31(1.95-2.73)	1.51(1.33-1.70)	1.35(1.16-1.55)
Both	No	Yes	2.28(1.90-2.71)	1.82(1.56-2.17)	1.46(1.28-1.65)
Both	No	No	2.30(1.87-2.78)	1.60(1.45-1.80)	1.44(1.30-1.61)

283

284

285

286 **Methods**

287 **Data availability.** All amino acid alignments and trees from this study are available as a

288 Supplementary Information file. GI numbers for sequences used in this study are

289 included in the taxon IDs.

290 **Code availability.** The code used in this analysis is included as a Supplementary
291 Information file. The authors place no restriction on its use.

292 **Data Collection.** BLASTP was performed against the NCBI NR Database using a series
293 of SQMO, OSC, SHC, and UbiH/COQ6 protein queries with an e-value cutoff of 10e-5
294 (Accession IDs: WP_027157849.1, NP_033296.1, AAH51055.1, WP_027156910.1,
295 NP_666118.1, WP_027157848.1, WP_027157848.1). The resulting hits were vetted for
296 conserved domains using PFAM²⁵. Because of the high overlap between OSC and SHC
297 BLAST hits, vetted results from the two searches were combined, sequences were
298 aligned with ClustalO²⁶, and a tree was made with RaxML²⁷, using an LG matrix and
299 100 rapid bootstraps. The results from this tree-building exercise clearly demarcated OSC
300 homologs from SHC, and we annotated the data accordingly.

301

302 Once the NCBI genes were vetted and annotated, subsampling was performed using
303 custom scripts. The script appends taxonomic information to each sequence from NCBI,
304 based on GI number, and then tabulates presence/absence data for all genes across all
305 taxa. Taxa are divided into clades based on a chosen taxonomic rank (in our scenario,
306 Order); if taxonomic data is missing from the NCBI taxonomy, the script automatically
307 looks one rank deeper. The script then determines whether any taxon in the clade
308 contains a copy of all gene queries. If multiple taxa contain all gene queries, the program
309 randomly selects one; otherwise, the program randomly selects a taxon with the highest
310 number of matching homologs. If the chosen taxon is missing a gene, but other species in
311 that clade have a copy of the gene, one sequence is randomly selected and added to the
312 dataset. For poorly represented clades in NCBI (amoebozoans, excavates, chlorophytes,

313 rhodophytes, rhizarians, and alveolates) we repeated this sampling at the one-per-genus
314 level. Still lacking rhizarian and alveolate data, we collected additional sequences from
315 the “SAR” clade using the Marine Microbial Eukaryotic Transcriptome Sequencing
316 Project (<http://marinemicroeukaryotes.org/>).

317

318 Following this process, certain taxa were removed from the dataset. Taxa were removed
319 if they contained multiple paralogs of a gene (e.g. most higher plants; 20% of all taxa), if
320 their higher-order taxonomy is contentious (e.g. haptophytes, glomeromycetes; 4% of all
321 taxa), or if they fell outside of their taxonomically accepted Class or Superphylum in one
322 or more analyses (e.g. nematodes, platyhelminthes; 6% of all taxa). We recovered SHC
323 proteins in many fungi, which suggests an interesting horizontal gene transfer event from
324 bacteria into eukaryotes. But because of poor resolution in the bacterial portion of the
325 SHC tree, we could not determine with confidence which lineage this transfer event
326 occurred from, and therefore chose to exclude fungal SHC from our study. Full datasets,
327 with the removed taxa indicated, are provided as Extended Data Figs. 1 & 2.

328

329 To analyze synteny, sequence (fasta) and annotation (gff3) files were downloaded from
330 NCBI for the following genome BioProjects: PRJNA82779, PRJNA242456,
331 PRJNA185587, PRJNA47603, PRJNA20997, PRJNA291650, PRJNA82927,
332 PRJNA63053, PRJDB3104, PRJNA21, PRJEA73721, PRJNA89087, PRJNA161599,
333 PRJNA203240, PRJNA19341, and PRJNA242472. The genomes were queried based on
334 accession numbers from the vetted SQMO and OSC proteins. Additional gene names
335 listed in Fig. 1B are based on the original annotation files associated with the genome.

336

337 **Protein alignment and tree building.** We ultimately created five curated datasets, which
338 were aligned with ClustalO: [1] SQMO (80 taxa, 841 characters), [2] SQMO +
339 UbiH/COQ6 (172 taxa, 1071 characters), [3] OSC (104 taxa, 1267 characters), [4] OSC +
340 SHC (174 taxa, 1475 characters), and [5] SQMO + OSC (116 taxa, 2166 characters).

341 Trees were built for [2] and [4] using RaxML on the CIPRES Science Gateway, with an
342 LG substitution model and 100 rapid bootstraps. Bayesian trees were constructed using
343 MrBayes²⁸ on the CIPRES Science Gateway, with MCMC sampling every 1000
344 generations for 4 million generations, or until convergence was reached according to the
345 *stopval* parameter (average standard deviation of split frequencies < 0.01). RaxML trees
346 for datasets [2] and [4] are provided in Extended Data Figs. 3 and 4 respectively, and
347 MrBayes trees are provided in Extended Data Figs. 5 and 6.

348

349 Molecular clocks were constructed using BEAST²⁹. We chose lognormal relaxed clocks,
350 using a yule process, an LG substitution model, and 4 gamma + invariant categories. In
351 dataset [5] we partitioned by gene, leaving substitution and clock models unlinked, but
352 trees linked. The clocks were calibrated using 18 fossils^{14,30-32} as described in Extended
353 Table 2. Calibration points were constrained as monophyletic, and three additional
354 monophyly constraints were set at nodes with poor resolution in our analyses, but high
355 resolution in previous multi-gene analyses: (1) Unikonta/Amorphea, (2) Bikonta
356 (including Bacterial Group II) and (3) non-oomycete stramenopiles. In analyses where
357 Excavata was treated as sister to Bikonta, we set an additional monophyly constraint
358 joining the two clades. BEAST MCMC chains were run for 20 million generations,

359 sampling every 1,000 generations. We re-ran all BEAST analyses a second time to test
360 for the reproducibility of our results, and then ran each analyses a third time sampling
361 only from the prior (i.e. ignoring molecular data), to verify that the dates we obtained
362 were not driven solely from fossil constraints. The results of these replications for
363 relevant nodes are plotted in Extended Data Fig. 7A. Finally, we generated calibration-
364 free molecular clocks using the RelTime method³³ to test the effect of fossil calibration
365 priors on our estimated divergence times. Plots of BEAST node dates mapped against
366 relative dates from RelTime (provided in Extended Data Fig. 7B) demonstrate a linear
367 relationship, suggesting no single fossil calibration is significantly altering the shape of
368 the tree.

369 **Methods References**

370

- 371 25 Finn, R. D. *et al.* The Pfam protein families database: towards a more sustainable
372 future. *Nucleic Acids Research* **44**, D279-D285 (2016).
- 373 26 Sievers, F. *et al.* Fast, scalable generation of high-quality protein multiple
374 sequence alignments using Clustal Omega. *Molecular Systems Biology* **7**, n/a-n/a,
375 doi:10.1038/msb.2011.75 (2011).
- 376 27 Stamatakis, A. RAxML-VI-HPC: maximum likelihood-based phylogenetic
377 analyses with thousands of taxa and mixed models. *Bioinformatics* **22**, 2688-2690
378 (2006).
- 379 28 Ronquist, F. & Huelsenbeck, J. P. MrBayes 3: Bayesian phylogenetic inference
380 under mixed models. *Bioinformatics* **19**, 1572-1574 (2003).
- 381 29 Drummond, A. J., Suchard, M. A., Xie, D. & Rambaut, A. Bayesian
382 Phylogenetics with BEAUti and the BEAST 1.7. *Molecular Biology and*
383 *Evolution* **29**, 1969-1973 (2012).
- 384 30 Yin, Z. *et al.* Sponge grade body fossil with cellular resolution dating 60 Myr
385 before the Cambrian. *Proceedings of the National Academy of Sciences* **112**,
386 E1453-E1460 (2015).
- 387 31 Benton, M. J. *et al.* Constraints on the timescale of animal evolutionary history.
388 *Palaeontologia Electronica* **18.1.1FC** (2015).
- 389 32 Gold, D. A., Runnegar, B., Gehling, J. G. & Jacobs, D. K. Ancestral state
390 reconstruction of ontogeny supports a bilaterian affinity for Dickinsonia.
391 *Evolution & Development* **17**, 315-324, doi:10.1111/ede.12168 (2015).
- 392 33 Battistuzzi, F. U., Billington-Ross, P., Murillo, O., Filipowski, A. & Kumar, S. A
393 Protocol for Diagnosing the Effect of Calibration Priors on Posterior Time

394 Estimates: A Case Study for the Cambrian Explosion of Animal Phyla. *Molecular*
395 *Biology and Evolution* **32**, 1907-1912 (2015).
396

397 **Extended Data Legends**

398 **Extended Data Figure 1:** Maximum likelihood (RAxML) tree, showing removal of
399 problematic SQMO sequences.

400 **Extended Data Figure 2:** Maximum likelihood (RAxML) tree, showing removal of
401 problematic OSC sequences.

402 **Extended Data Figure 3:** Maximum likelihood (RAxML) tree from vetted SQMO
403 dataset.

404 **Extended Data Figure 4:** Maximum likelihood (RAxML) tree from vetted OSC dataset

405 **Extended Data Figure 5:** Bayesian (MrBayes) tree from vetted SQMO dataset.

406 **Extended Data Figure 6:** Bayesian (MrBayes) tree from vetted OSC dataset.

407 **Extended Data Figure 7:** Reproducibility of BEAST runs, and relationship between
408 BEAST and RelTime trees.

409 **Extended Table 1:** Distribution of marginal probabilities for all molecular clock analyses
410 (as percentages), binned by geologic time.

411 **Extended Table 2:** Fossil calibration points used in molecular clock. An asterisk (*)
412 denotes a calibration point only used in the UbiH/Coq6 outgroup. A caret (^) signifies
413 that monophyly was not enforced on this clade as a prior in Bayesian analysis.

414

415





