



1 **Abstract**

2

3 The supergroup Amoebozoa unites a wide diversity of amoeboid organisms and  
4 encompasses enigmatic lineages recalcitrant to modern phylogenetics. Deep divergences,  
5 taxonomic placement of some key taxa and character evolution in the group largely  
6 remain poorly elucidated or controversial. We surveyed available Amoebozoa genomes  
7 and transcriptomes to mine conserved putative single copy genes, which were used to  
8 enrich gene sampling and generate the largest supermatrix (824 genes) in the group to  
9 date. We recovered a well-resolved and supported tree of Amoebozoa, revealing novel  
10 deep level relationships and resolving placement of enigmatic lineages congruent with  
11 morphological data. In our analysis the deepest branching group is Tubulinea. A recent  
12 proposed major clade Tevosa, uniting Evosea and Tubulinea, is not supported. Based on  
13 the new phylogenetic tree, paleoecological and paleontological data as well as data on the  
14 biology of presently living amoebozoans, we hypothesize that the evolution of  
15 Amoebozoa probably was driven with the need to disrupt and graze on microbial mats -  
16 a dominant ecosystem of the mid-Proterozoic period of the Earth history.

17

18 **Keywords:** Amoebozoa, phylogenomics, flagellum, eukaryotes, genome, transcriptome

## 1 Introduction

2  
3 The supergroup Amoebozoa<sup>1</sup> comprises a variety of amoeboid lineages; namely,  
4 naked lobose amoebae (which are “archetypal” amoebae), testate lobose amoebae,  
5 mycetozoa, anaerobic archamoebians and a heterogeneous assemblage of flattened  
6 amoeboid, branching reticulate or flagellated organisms; presently known as Variosea.  
7 Amoebozoa holds a key evolutionary position, being the closest known relative of  
8 Obazoa that, among other organisms, includes humans<sup>2,3</sup>. Resolving the phylogenetic  
9 tree of this lineage is critical for answering important questions pertaining to the  
10 evolutionary origin of Amoebozoa, as well as for further clarification of the root of the  
11 eukaryotic tree<sup>3-8</sup>.

12  
13 Our understanding of the evolution and taxonomy of amoeboid protist originally  
14 conceived from cytological, morphological and life cycle evidence<sup>9,10</sup>. Early studies  
15 based on small subunit rDNA (18S) gene indicated the polyphyly of naked amoebae  
16 (gymnamoebae) and formed the basis of our understanding of the supergroup  
17 Amoebozoa<sup>1,11,12</sup>. The assemblage of Amoebozoa grew in membership, albeit with little  
18 improved resolution; or sometimes with conflicting hypotheses pertaining to within-  
19 group relationships (e.g.,<sup>13-19</sup>). This led to subsequent revisions and reevaluation in  
20 attempts to combine morphological and molecular characters and find synapomorphic  
21 characters of major clades<sup>20-23</sup>. While this achieved major progress in our overall  
22 understanding of the group, much of the deep and intermediate relationships and  
23 placement of some groups of uncertain phylogenetic affinities (so-called *incertae sedis*  
24 taxa) remained elusive. Multigene studies, varying in breadth and depth of gene and  
25 taxon sampling, managed to overcome many of the challenges of single-gene  
26 reconstructions; and they resolved some of the long-standing evolutionary questions in  
27 the group<sup>4,24-27</sup>. A recent phylogenomic study by Kang et al.<sup>4</sup> reported a deep level  
28 phylogeny of Amoebozoa based on large taxon sampling. However, the placements of  
29 some *incertae sedis* lineages were not entirely resolved. For some groups, other  
30 phylogenomic studies reported conflicting relationships<sup>25,26,28</sup>.

31  
32 The conflict in existing phylogenomic studies can be attributed partially to  
33 limitations of taxon and gene sampling as well as the methodology. Kang et al.<sup>4</sup> used  
34 large taxon sampling, but included only a small fraction of data (325 genes), from the  
35 vast amount of transcriptomic and genomic data available, based on commonly used  
36 genetic markers in eukaryotes. There are data suggesting that taxon sampling alone is not  
37 sufficient to resolve deep divergences in ancient lineages that might have undergone  
38 rapid radiations<sup>29</sup>. The age of Amoebozoa is estimated to be over a billion years, and the  
39 probable origin of the group is dated back to the mid-Proterozoic period<sup>30,31</sup>. Therefore,  
40 in order to infer deep evolutionary divergences not only increased taxon sampling, but  
41 also more representative genetic sampling along with the application of appropriate  
42 models and methods are essential.

43  
44 In this study, we sampled putative single copy gene markers from genome-wide  
45 assays, increased taxon sampling and produced the largest amoebozoan supermatrix to  
46 date. This large dataset enabled us to recover a well-resolved and supported tree of the

1 Amoebozoa. In addition, we uncover a well-corroborated novel deep-level relationship  
2 and resolved the placement of some *incertae sedis* lineages.

## 3 4 5 **Results**

### 6 7 **The Tree of Amoebozoa**

8  
9 We recovered a monophyletic tree of Amoebozoa that is well resolved and  
10 supported in every one of our analyses (Figs. 1, S1-S4). Our datasets, with and without  
11 fast-evolving sites removed (analyzed using the complex model in IQ-TREE) recovered  
12 all well-established major subclades of Amoebozoa including Discosea, Archamoebae,  
13 Cutosea, Eumycetozoa, Variosea and Tubulinea with full support (Figs. 1, S1). The two  
14 well-known long-branch lineages, Archamoebae and Cutosea, were placed in their  
15 respective correct phylogenetic positions without removal of fast evolving sites in our full  
16 dataset (Fig. 1). Removal of fast evolving, rate categories, in IQ-TREE neither affected  
17 the topology nor improved support values (Fig. S1). In the RAxML analysis, the accurate  
18 placement of Archamoebae and Cutosea, required removal of six fast evolving rate  
19 categories (38%) from the full dataset (Fig. S2); but resulted in the same final tree  
20 configuration. The RAxML tree had generally lower supported branches but was  
21 congruent with the topology of the trees inferred using IQ-TREE (Figs. 1, S1, S2). A  
22 similar reduced dataset was analyzed using Bayesian inference, which yielded similar  
23 topology despite lack of convergence in our PhyloBayes analysis (data not shown). Kang  
24 et al. <sup>4</sup> also reported similar topologies among their ML and PhyloBayes trees despite  
25 limited number of chains used and lack of convergence in some of their PhyloBayes  
26 analyses. Due to the high computational demand, Bayesian inference was not feasible  
27 with our large dataset. The consistency of tree topologies across methods and algorithms  
28 used, as well as the placement of long-branch taxa (Archamoebae and Cutosea) without  
29 removal of fast evolving sites in IQ-TREE (likely due to complex model used),  
30 demonstrates the robustness of our result.

31  
32 In our phylogenomic tree, all major clades are congruent with previous published  
33 topologies <sup>4,24-26</sup>. Moreover, our phylogenomic tree has well-corroborated relationships;  
34 and the recovery and placement of enigmatic taxa are more stable (Figs. 1, S1, S2). Our  
35 results yielded improved support for the Flabellinia and Thecamoebida clades compared  
36 to a previous comparable phylogenomic study <sup>4</sup>. We have recovered for the first time a  
37 fully supported monophyletic clade encompassing two *incertae sedis* taxa, *Vermistella*  
38 and *Stygamoeba*. Both these lineages were placed in the order Stygamoebida based on  
39 morphological evidence <sup>22</sup>. The monophyly and placement of this order in the tree of  
40 Amoebozoa has not been resolved in previous multigene analyses (e.g., <sup>4</sup>). In our tree  
41 Stygamoebida clade forms a sister group relationship with Thecamoebida with full  
42 support (Fig. 1). We also find some discrepancies between our tree (Fig. 1) and that of  
43 Lahr et al. <sup>5</sup> in the branching order of the Tubulinea clade, albeit with similar taxon  
44 sampling for this clade. Our analysis shows clade Corycida as the most basal Tubulinea  
45 lineage similar to that of Kang et al. <sup>4</sup> phylogeny (Fig. 1). *Nolandella* sp., a member of

1 Euamoebida, did not group with *Amoeba proteus* and *Copromyxa protea* in our analysis,  
 2 but formed an independent lineage (Fig. 1).  
 3  
 4



5 **Figure 1.** Genome wide phylogeny of the Amoebozoa inferred using Maximum  
 6 likelihood (ML) in IQ-TREE with LG+G4+C60+F model of evolution. The data matrix  
 7 used to infer this tree consisted of 113,910 amino acid sites from the full dataset, derived  
 8 from 824 genes and 113 taxa including 10 outgroup taxa. Clade supports at nodes are ML  
 9 IQ-TREE 1000 ultrafast bootstrap values obtained using the same model. All branches  
 10 are drawn to scale except a branch leading to Archamoebae, and *Sapocribum*  
 11 *chincoteaguense* and *Parvamoeba monoura*, that were reduced to one-third and half,  
 12 respectively.  
 13

14 **A Novel Deep Split of the Amoebozoa**

15  
 16 Our analysis for the first time revealed a novel, well-supported deep split of  
 17 Amoebozoa; not reported in previous phylogenomic studies. Amoebozoa is split into two  
 18 fully supported major subclades: Tubulinea and a second one comprised of the remaining  
 19

1 major subclades including Evosea (Eumycetozoa, Variosea, Archamoebae, and Cutosea)  
2 and Discosea (Figs. 1, S1, S2). This branching is different from a finding in a recent  
3 phylogenomic study that reported a split between Discosea and Tevosa  
4 (Evosea+Tubulinea)<sup>4</sup>. Tevosa is not supported in our analyses, including analyses with  
5 removal of fast sites. On the other hand, the deep split (Evosea+Discosea vs. Tubulinea)  
6 observed in our phylogenomic tree is supported in all analyses of our data sets. The deep  
7 split receives almost full support in our internode certainty (IC) analyses as implemented  
8 in QuartetScores (1.00) and RAxML (0.979) (Figs. S3, S4). AU test of our topology,  
9 comparing alternative topologies with Tevosa and a traditional deep relationship uniting  
10 Discosea and Tubulinea (Lobosa), showed that the newly recovered deep split has the  
11 highest p-value (p-AU = 0.947). Hypothesis Lobosa was rejected (p-AU = 0.000278),  
12 while Tevosa cannot be rejected with p-value just above threshold (p-AU = 0.0564). For  
13 convenience, we suggest a new name for the deep split (Discosea+Evosea) clade; i.e.,  
14 Divosa, a term derived from a combination of the name of the two clades.

15

16

## 17 **Discussion**

18

### 19 **Targeted Genome-Wide Data Enrichment for Phylogenomics of Amoebozoa**

20

21 Despite the large number of RNA-Seq data generated in recent studies<sup>4,24-26</sup>, only  
22 a small fraction of this data has been utilized in phylogenomic analyses. To increase it,  
23 we compiled a total of 1559 markers using genome-derived protein coding genes from  
24 113 amoebozoan genomes and transcriptomes. Using putative single copy markers,  
25 primarily derived from Amoebozoa genomes, has enabled us to introduce highly  
26 conserved markers with phylogenetic signal corroborating morphology- and  
27 phylogenomic-based amoebozoan hypotheses<sup>4,24</sup>. While single-copy genes identified in  
28 some genomes might not always apply to others, a previous phylogenomic study with  
29 seed plants, based on single copy markers resulted in more resolved phylogeny both at  
30 shallow and deep nodes<sup>32</sup>. In this study, we followed a stringent approach aided by  
31 automated and manual curation of markers, selected from the above-mentioned dataset to  
32 build the largest supermatrix (823 genes) in the Amoebozoa. With this approach, we  
33 substantially increased the total number of genes used in Amoebozoa phylogenomics.  
34 Our analysis yielded consistent and well-corroborated topologies, despite whether we  
35 included or excluded fast evolving sites (Figs. 1, S2). The robustness of our phylogeny is  
36 also corroborated with the high support values from internode certainty analysis (Figs.  
37 S3, S4). One of the evident results of this approach is the first time phylogenomic  
38 recovery of the monophyly of the taxon Stygamoebida, earlier supported only at the  
39 morphological level<sup>22,23</sup> and a recovery of a novel deep split divergence of Amoebozoa.

40

### 41 **Unraveling deep divergence of Amoebozoa**

42

43 A recent phylogenomic study by Kang et al.<sup>4</sup>, though based on a slightly smaller  
44 taxon sampling, proposed a split of the Amoebozoa supergroup into two major subclades:  
45 Tevosa (Evosea+Tubulinea) and Discosea. By contrast, in our study Evosea robustly  
46 groups as sister clade to Discosea (Figs. 1, S1, S2). Both phylogenetic hypotheses,

1 ‘Tevosa’ and Divosa, receive high statistical support in their and our study, respectively  
2 (see Fig. 1, <sup>4</sup>). In phylogenomic analyses, it is common to see that short subtending deep  
3 nodes receive high statistical support <sup>33</sup>. Amoebozoan deep nodes are characterized by  
4 very short branch lengths, an indication of limited supporting characters, or possible  
5 ancient rapid diversification. Strong statistical support at these levels of nodes does not  
6 necessarily mean that the inferred relationships are correct. Statistical indices such as  
7 bootstrap values and Bayesian posterior probabilities only assess sampling effects, and  
8 give an indication of tree reliability that is dependent on the data and the method <sup>34</sup>. This  
9 can partially explain why these short-branch, deep nodes in Amoebozoa phylogenomic  
10 studies tend to collapse, or vary, depending on the method of analysis or the composition  
11 of the gene/taxon sampling <sup>4,24-26</sup>. Certainly, caution still must be taken when interpreting  
12 ancient divergences, because results can be muddied by noise (e.g., gene history <sup>35</sup> or lack  
13 of signal due to rapid radiation <sup>29</sup>). However, the support of the split recovered in the  
14 present study is high and originates from different lines of evidence.

15  
16 It is possible to note that in many lineages trophozoites of Discosea and Variosea  
17 are more similar to each other rather than to Tubulinea. Certainly, the morphology of  
18 presently living amoeboid organisms is derived and adaptive, but generally it is possible  
19 to say that members of Divosa lineage share more morphological similarity between each  
20 other rather than with the Tubulinea lineage. For example, amoebae of the genus  
21 *Flamella*, belonging to the class Variosea, by their morphology may be easily confused  
22 with some discosean amoebae (e.g., <sup>36</sup>); the same is true for individual trophozoites of  
23 many mycetozoa species, showing flattened body shape and pointed subpseudopodia  
24 <sup>37,38</sup>. Cells of amoebae belonging to the genus *Squamamoeba* (the taxon of Cutosea),  
25 sometimes resemble *Korotnevella* (Discosea) in their overall morphology, hence, being  
26 differently organized at the cytological level <sup>39</sup>. At the same time, none of discosean or  
27 variosean lineages show the morphology resembling that of, e.g. Amoebida, or alteration  
28 of the locomotive morphology from flattened to tubular, which is a general characteristic  
29 of Tubulinea <sup>20,22</sup>. To certain extent, the return to the tubular body shape, subcylindrical  
30 in cross-section occurs among amoeboid representatives of Archamoebae; however, this  
31 might be mostly related with their specific lifestyle (parasites or pelobionts). In addition  
32 the pattern of pseudopod formation (e.g., the tendency to show eruption of the hyaline  
33 cytoplasm in the frontal area of the cell) makes them to be significantly different from  
34 that in Tubulinea (see <sup>40</sup>).

### 35 **Mid-Proterozoic environment – the driving force for the origin of Amoebozoa**

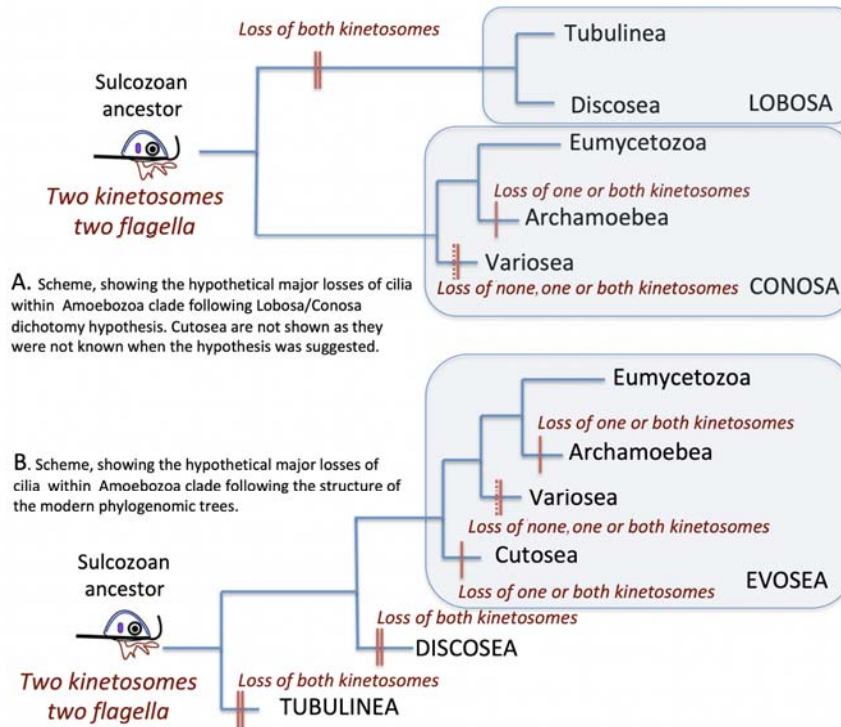
36 The flagellum (cilium) is a highly conserved complex structure that is believed to  
37 have originated only once, and be ancestral to all eukaryotes <sup>2,41,42</sup>. Amoebozoa are  
38 remarkable in that the two basal phylogenetic lineages, Tubulinea and Discosea, have  
39 entirely lost cilia, kinetosomes (basal bodies) and associated root structures; while a  
40 derived major clade, Evosea, contains a handful of ciliated lineages in a few branches  
41 intermingled among amoeboid lineages <sup>21,22</sup>. The loss of cilia and associated structures in  
42 the majority members of Amoebozoa is one of the biggest mysteries pertaining to their  
43 origin and evolution.

1           In ciliated members of Amoebozoa, the ciliary apparatus is characterized by a  
2 specific arrangement of root structures, which includes an incomplete (Variosea and  
3 Mycetozoa) or complete (Archamoebae) cone of microtubules extending from the  
4 kinetosome to the nucleus<sup>43</sup>. In early interpretations, this conical arrangement of  
5 microtubules was considered to be homologous to the ciliary root system of  
6 Opisthokonta; which, together with other morphological and molecular evidence, gave  
7 rise to the “Unikonta” hypothesis<sup>2,44,45</sup>. In this model, the hypothetical ancestor of  
8 Amoebozoa was considered to be an organism with a single emergent cilium, resembling  
9 *Phalansterium* or *Mastigamoeba* in cellular organization<sup>46,47</sup>. This lineage, combining  
10 Amoebozoa and Opisthokonta, has been proposed as an alternative to that of the bikonts,  
11 with two emerging cilia; which included the rest of the eukaryotic groups. Cavalier-Smith  
12<sup>2</sup> argued that among unikonts, paired kinetosomes (when present) resulted from  
13 convergent evolution rather than common ancestry with bikonts. Molecular and  
14 morphological analyses provided certain indications that the microtubular structures in  
15 Amoebozoa, and Opisthokonta may not be homologues<sup>43,48</sup>. However, further  
16 development of molecular phylogeny provided evidence for the basal position of bikont  
17 organisms in the tree of eukaryotes<sup>3,49,50</sup>. Thereafter, the general consensus nowadays is  
18 that hypothetical common ancestor of Amoebozoa, was a bikont organism<sup>43,51,52</sup>.  
19 Several authors (e.g.,<sup>3,43,49,50</sup>) hypothesised that the presumable common ancestor was a  
20 ventrally grooved biciliate gliding flagellate, capable of producing filose ventral  
21 pseudopodia and possessing a relatively complex organization of the cell. That is, a cell  
22 possessing two cilia with kinetosomes and root structures, ventral groove supported with  
23 microtubules and dorsal pellicle – the so called “sulcozoan ancestor”. Its name originates  
24 from Sulcozoa – a phylum of protists established by Cavalier-Smith<sup>43</sup> that combines a  
25 heterogenous assemblage of early evolving eukaryotic lineages. Cavalier-Smith  
26 suggested that “opisthokonts and Amoebozoa evolved from sulcozoan ancestors by two  
27 independent losses of the pellicular dense layers and of the ventral groove, which in both  
28 cases would allow pseudopods to develop anywhere on the cell surface” (op. cit.).

29           The origin and further evolution of Amoebozoa in this hypothesis presumes the  
30 loss of both cilia and kinetosomes in Lobosa (Tubulinea and Discosea) and of the  
31 posterior cilium and one kinetosome in most of the ancestors of Conosa - Archamoebae,  
32 Variosea and Eumycetozoa; Cutosea were not known at that time (e.g.,<sup>3,49,50</sup>). This  
33 evolutionary scenario was rather logical and is illustrated in Figure 2A. However, the  
34 Lobosa/Conosa dichotomy was doubted based on some 18S gene phylogenies<sup>27</sup>; and it  
35 subsequently failed to garner support in wide-scale phylogenomic studies<sup>4,24,25</sup>, as well  
36 as in the present study. This makes the model of multiple losses more complicated,  
37 because under the new tree configuration, we have to suggest subsequent partial or  
38 complete loss of cilia and related structures in all but one branch of Amoebozoa. This  
39 hypothetical scenario is illustrated in Figure 2B. It remains unclear why the hypothesized  
40 ancestor of Amoebozoa, being initially a quite complex biciliated organism, underwent  
41 such a massive loss (or substantial simplification) of cilia-related structures in almost all  
42 evolutionary lineages of Amoebozoa, and what was the driving force for such a  
43 reduction.

44





1

2 **Figure 2.** A scheme illustrating the loss of kinetosomes and cilia under the different  
3 evolutionary hypotheses (A and B). Vertical hash marks on branches show loss of  
4 kinetosomes (the number lost as designated by labels on the diagram) depending on the  
5 lineage.

6 Several studies based on molecular dating analysis correspondingly placed the  
7 origin of Amoebozoa to the Mesoproterozoic period, which means 1250 – 1624 mya<sup>31,53</sup>.  
8 It means that the early evolution of Amoebozoa took place at the period when the  
9 biosphere was dominated with microbial biofilms – sheets of bacteria, embedded in  
10 extracellular polymeric substances, covering almost every possible substrate<sup>54</sup>. Being  
11 initially rather simple, biofilms further evolved in complex microbial mats, comprising  
12 different prokaryotic organisms, showing concerted activities and intimate interactions  
13 between various microbial metabolisms<sup>55</sup>. The oldest mats are dated to approximately  
14 3.5 billion years ago, and the noonday of mats covers the mid-Proterozoic period<sup>56,57</sup>,  
15 which roughly corresponds to the estimate of the potential age of Amoebozoa.

16 Formation of a microbial biofilm, among other structural and biogeochemical  
17 features, can be explained as an adaptation that increases survival of bacteria to avoid  
18 predation<sup>58,59</sup>. The probable size of the bacterivorous biflagellate ancestor of Amoebozoa  
19 was relatively small, likely no larger than that of the existing representatives of the  
20 CRuMs clade (e.g., *Mantamonas*) or ‘Excavates’ (metamonads or *Malawimonas*), which  
21 is within the general size range of 2-20  $\mu\text{m}$ . These organisms were able to phagocytize  
22 solitary bacteria, but consumption of microorganisms embedded in an intact microbial  
23 mat probably was beyond their capacity, as well as this is beyond the capacity of the  
24 modern flagellates of comparable size<sup>60,61</sup>. Feeding on bacteria, major constituents of the

1 microbial mats (the dominant food source in the mid-Proterozoic environment), required  
2 increment in the body size and acquisition of special adaptations allowing them to ingest  
3 filamentous food. However, the latter was again related to the body size, because the  
4 filament, even compacted in some way, must be ingested – i.e., appear inside the cell.

5         Due to Reynolds number limitation<sup>62,63</sup>, the increment in the body size makes  
6 ciliary motility less adaptive due to loss of efficiency. Thus, from an adaptive aspect, an  
7 amoeboid lifestyle might be a way to increase the body size while retaining a motility  
8 function, no longer dependent on cilia. An amoeboid organization also could gain the  
9 adaptive capacity to disrupt microbial mats and graze, feeding on bacteria within the  
10 mats. This adaptation would provide access to the dominant food source in the biosphere  
11 of the mid-proterozoic eon. Indeed, presently, naked amoebae are known as one of the  
12 primary grazers of bacterial biofilms<sup>64-66</sup>. Moreover, they not only just graze and  
13 phagocytize prey in the mats, but also disrupt them, making their content available for  
14 other organisms<sup>67,68</sup>. Finally, in addition to the advantage of feeding on bacterial mats  
15<sup>69,70</sup>, it is also possible that an increase in body size alleviated pressure of predation by  
16 other organisms on the last Amoebozoan common ancestor (LACA), which for some  
17 time provided it an adaptive advantage and allowed rapid proliferation and differentiation  
18 of Amoebozoa in the mid-Proterozoic environment.

19         Hence, we hypothesise that the adaptive value of amoeboid locomotion and  
20 concomitant grazing potential on the dominant food source in the mid-proterozoic  
21 biosphere – the microbial mats – favoured the evolution of the Amoebozoa. They  
22 probably successfully solved this task by the increment of body size. However, at the  
23 same time, the efficiency of flagellar locomotion was highly reduced or lost; and this  
24 resulted in the multiple suspensions of the flagellar apparatus, which is completely absent  
25 in two major current amoebozoan lineages – Tubulinea and Discosea (Fig. 2). The  
26 modern configuration of the Amoebozoan tree, which rejects the Lobosa/Conosa  
27 dichotomy and suggests a subsequent branching of lineages (with either Tubulinea or  
28 Discosea at the base), leaves open a major question. That is, was the last Amoebozoa  
29 common ancestor an amoeboflagellate, with the domination of amoeboid movement  
30 based on the microtubular cytoskeleton; or was the flagellum-related structures and  
31 microtubular locomotive system entirely suppressed? If the latter case is true, then it  
32 probably drove the ancestral amoebozoan to switch to the acto-myosin movement, as  
33 found in modern representatives of naked and testate lobose amoebae. Probably, the  
34 answer to this question may be obtained by the analysis of gene content and the level of  
35 flagellum-related gene expression in the amoebozoan genomes. However, the dataset  
36 available for quality analysis remains limited in this group of protists and requires further  
37 accumulation prior to conclusive study.

38

## 39 **Methods**

40

### 41 **Transcriptome Assembly and Contamination Examination**

42

43         All transcriptome data used in this study were assembled using a bioinformatics  
44 pipeline described in Tekle and Wood<sup>25</sup>. As a precautionary measure for contamination,

1 high-quality data generated from single cell or monoclonal cultures, and without history  
2 of contamination, were prioritized in our data collection. We also checked highly  
3 conserved genes (e.g., small subunit rDNA and cytoskeletal genes) for assembled  
4 transcriptomes to check the identity of the species. Species suspected to have been  
5 contaminated (e.g., *Ripella* sp. DP13-Kostka) or with low- or poor-quality transcriptome  
6 data (see below) have been removed from the final analysis. Assembled contigs were  
7 translated into protein sequences using TransDecoder  
8 (<https://github.com/TransDecoder/TransDecoder/wiki>).  
9

## 10 **Taxon and gene sampling**

11  
12 A total of 107 amoebozoans representing the vast diversity of the supergroup and  
13 10 outgroup taxa from a closely related clade, Obazoa, were included in our initial  
14 analysis (Table S1). Four ingroup taxa including *Parvamoeba rugata*, *Centropyxis*  
15 *aculeata*, *Hyalosphenia elegans* and *Grellamoeba robusta*, were removed from the final  
16 dataset due to poor data quality. A recent phylogenomic study<sup>5</sup> that focused on testate  
17 amoebae (clade Tubulinea) reported a topology of Tubulinea that differed from that of  
18 Kang et al.<sup>4</sup>. To explore these discrepancies further, and assess the impact of taxon  
19 sampling on branching order of Tubulinea clade and its position within the Amoebozoa  
20 phylogeny, we added more slowly evolving taxa to Tubulinea. The final supermatrix  
21 consisted of 113 taxa including the outgroup taxa (Table S1).  
22

23 A genome wide gene sampling approach using available amoebozoan genomes  
24 was employed to identify single copy markers. Previous phylogenomic studies have used  
25 conserved phylogenetic markers commonly found in a wide range of eukaryotic diversity  
26<sup>4,24</sup>. In this study we used a series of bioinformatics steps to maximize gene sampling in  
27 the Amoebozoa. We conducted a whole genome comparison of three well-annotated  
28 amoebozoan genomes, *Acanthamoeba castellanii*, *Dictyostelium discoideum* and  
29 *Entamoeba histolytica*, to extract commonly shared protein-coding genes among these  
30 genomes in OrthoVenn<sup>71</sup>. Inclusion of *E. histolytica* greatly reduced the number of  
31 shared genes by 40% because this amitochondriate parasitic species has a comparably  
32 much reduced genome to the free-living amoebae. For this reason, to be more  
33 representative, further comparative analysis was done using *A. castellanii* and *D.*  
34 *discoideum* as reference genomes to mine single-copy genes. Using this approach, we  
35 identified 1559 putative single copy genes that were used as a query to search  
36 orthologous genes from ingroup and outgroup taxa.  
37

38 We used NCBI-BLAST with e-value threshold of  $10^{-15}$  to retrieve homologous  
39 sequences from transcriptomes or genomes of our selected taxa. From this analysis,  
40 sequences with best e-value scores were retained for each taxon. The retained sequences,  
41 for each taxon and gene, were compiled and aligned using a sequence alignment tool,  
42 MAFFT, with default setting<sup>72</sup>. These alignments were then trimmed in TrimAl<sup>73</sup> using  
43 “automated1” setting provided by the program. To inspect potential paralogs from each  
44 gene, we inferred single gene trees using IQ-TREE with the best-fit model automatically  
45 fast selected by ModelFinder<sup>74</sup>. Both single gene trees and their corresponding  
46 alignments were then inspected manually for paralogy and other anomalies related to

1 alignment accuracy, sequence length and fast evolving lineages ((Single gene alignment  
2 and trees available for review on this link: <https://www.dropbox.com/>). We applied strict  
3 gene selection criteria that included removal of anomalous grouping (e.g., lineages that  
4 grouped with outgroup or wrong (unexpected) phylogenetic position with >90%  
5 bootstrap support) and genes that showed paralogy (duplication) signs. To mitigate the  
6 impact of long-branch attraction during phylogenetic reconstruction, we removed genes  
7 that contained three or more long-branch lineages. Two exceptions for this approach were  
8 the well-known long-branch lineages, Cutosea and *Entamoeba*, that were kept in all of  
9 our analyses. These two lineages were retained since all their representatives are mostly  
10 long-branches. They are also indirect indicators of noise in a data matrix since their  
11 correct placement usually requires removal of fast-evolving sites due to the effect of  
12 long-branch attraction. Following these criteria, we retained a total of 824 gene clusters  
13 in the final dataset. Orthologous group numbers were assigned for each gene cluster using  
14 ublast in USEARCH <sup>75</sup> with e-value  $10^{-10}$ . We used the OrthoMCL database to generated  
15 ortholog group numbers <sup>76</sup> (Table S2).

16

## 17 **Supermatrix Construction and Tree Inference**

18

19 The alignments from 824 genes were concatenated into an initial supermatrix  
20 containing 198,280 amino acid sites and 117 taxa using a customized R script. Taxa with  
21 over 80% gappy sites were removed, which resulted in exclusion of 4 lineages  
22 (*Parvamoeba rugata*, *Centropyxis aculeata*, *Hyalosphenia elegans*, *Grellamoeba*  
23 *robusta*). Constant sites, and sites with more than 50% missing data, were removed from  
24 this alignment, and the resulting supermatrix retained 113,910 amino acid sites and 113  
25 taxa for the full dataset.

26

27 Phylogenomic analyses of the final datasets were conducted in IQ-TREE – an  
28 efficient tool to analyze large datasets by the maximum likelihood (ML) method <sup>74</sup>. All  
29 IQ-TREE analyses were performed using LG+G4+C60+F model, with 1000 replicates  
30 for ultrafast bootstrap, which allowed full profile mixture model C60 and Gamma rate  
31 heterogeneity across sites. We also analyzed our dataset in RAxML v.8.2.X <sup>77</sup> using  
32 PROTGAMMALG4X model; branch support was estimated from 1000 rapid bootstrap  
33 pseudoreplicates.

34

35 Fast-evolving sites and taxa are known to be problematic for tree inference due to  
36 saturation of substitutions and subsequent convergent evolution resulting in long-branch  
37 attraction (LBA) and other systematic errors. To test the effects of these types of errors  
38 on our phylogenomic analysis, we performed a site removal assay in which each site of  
39 the supermatrix was assigned to one of 16 categories based on its rate from IQ-TREE.  
40 This was performed using a posterior mean site frequency (PMSF) model with mixture  
41 model C60 and 16 discrete rate categories of sites. For this analysis, we used the tree  
42 from full dataset inferred above as a guide tree. The impact of fast evolving sites on  
43 resulting phylogenies was assessed by subsequent removal of fast categories of sites (up  
44 to 6 categories). In IQ-TREE our full dataset was analyzed with 3 categories removed  
45 using PMSF model with a guide tree inferred from the complex model (LG+G4+C60+F)

1 mentioned above. In RAxML, 3 and 6 fast site categories were removed and analyzed  
2 using the same model as above.

### 4 Internode Certainty Analysis and Hypothesis Testing

6 As alternative to bootstrap branch support from IQ-TREE, we calculated  
7 internode certainty (IC) scores using the program QuartetScores<sup>78</sup>. This approach  
8 calculated IC scores from the frequencies of quartets, which can correct for the missing  
9 taxa using a set of trees. For this analysis, we used 1000 bootstrap trees generated from  
10 LG+G4+C60+F model in IQ-TREE with our full dataset. Alternatively, we used RAxML  
11 to estimate the degree of certainty for internodes and tree topology for bipartitions with  
12 PROTGAMMALG4X model<sup>79</sup>.

14 We used Approximately Unbiased (AU) tests<sup>80</sup> to test alternate tree topologies  
15 pertaining to the deep node hypotheses Divosa (this study), Tevosa (Kang et al. 2017)  
16 and Lobosa<sup>27</sup> with the full dataset (113,910 sites). Two loosely constrained topologies  
17 Tevosa ([Tubulinea+Evosea]+Discosea) and Lobosa ([Discosea+Tubulinea]+Evosea)  
18 were optimized under LG+G4+F+C60 in IQ-TREE. These optimized trees were  
19 compared with our tree (Divosa, ([Discosea+Evosea], Tubulinea) using AU test with  
20 10,000 RELL bootstrap replicates<sup>81</sup>. The hypotheses that had  $p\text{-AU} \geq 0.05$  within the  
21 95% confidence interval could not be rejected.

### 24 References

- 26 1 Cavalier-Smith, T. A revised six-kingdom system of life. *Biological Reviews of*  
27 *the Cambridge Philosophical Society* **73**, 203-266 (1998).
- 28 2 Cavalier-Smith, T. The phagotrophic origin of eukaryotes and phylogenetic  
29 classification of protozoa. *International Journal of Systematic and Evolutionary*  
30 *Microbiology* **52**, 297-354 (2002).
- 31 3 Brown, M. W. *et al.* Phylogenomics demonstrates that breviate flagellates are  
32 related to opisthokonts and apusomonads. *Proc Biol Sci* **280**, 20131755,  
33 doi:10.1098/rspb.2013.1755 (2013).
- 34 4 Kang, S. *et al.* Between a Pod and a Hard Test: The Deep Evolution of Amoebae.  
35 *Mol Biol Evol* **34**, 2258-2270, doi:10.1093/molbev/msx162 (2017).
- 36 5 Lahr, D. J. G. *et al.* Phylogenomics and Morphological Reconstruction of  
37 Arcellinida Testate Amoebae Highlight Diversity of Microbial Eukaryotes in the  
38 Neoproterozoic. *Curr Biol* **29**, 991-1001 e1003, doi:10.1016/j.cub.2019.01.078  
39 (2019).
- 40 6 Yoon, H. S. *et al.* Broadly sampled multigene trees of eukaryotes. *BMC*  
41 *Evolutionary Biology* **8**, 14 (2008).
- 42 7 Burki, F. *et al.* Phylogenomics reshuffles the eukaryotic supergroups. *PLoS ONE*  
43 **2**, e790 (2007).
- 44 8 Parfrey, L. W. *et al.* Broadly sampled multigene analyses yield a well-resolved  
45 eukaryotic tree of life. *Systematic biology* **59**, 518-533, doi:syq037 [pii]  
46 10.1093/sysbio/syq037 (2010).

- 1 9 Page, F. C. The classification of 'naked' amoebae (Phylum Rhizopoda). *Arch.*  
2 *Protistenkd.* **133**, 199–217 (1987).
- 3 10 Rogerson, A. & Patterson, D. J. The Naked Ramicristate Amoebae  
4 (Gymnamoebae). In: Lee, J.J., Leedale, G.F., Bradbury, P. (Eds.), *An Illustrated*  
5 *Guide to the Protozoa, 2nd ed. Society of Protozoologists, Lawrence, Kansas*, pp.  
6 1023–1053 (2002).
- 7 11 Amaral Zettler, L. A. *et al.* Microbiology: Eukaryotic diversity in Spain's River of  
8 Fire. *Nature* **417**, 137 (2002).
- 9 12 Cavalier-Smith, T. & Chao, E. E. Molecular phylogeny of the free-living  
10 archezoan *Trepomonas agilis* and the nature of the first eukaryote. *Journal Of*  
11 *Molecular Evolution* **43**, 551-562 (1996).
- 12 13 Tekle, Y. I. *et al.* Phylogenetic placement of diverse amoebae inferred from  
13 multigene analyses and assessment of clade stability within 'Amoebozoa' upon  
14 removal of varying rate classes of SSU-rDNA. *Molecular phylogenetics and*  
15 *evolution* **47**, 339-352 (2008).
- 16 14 Fahrni, J. F. *et al.* Phylogeny of lobose amoebae based on actin and small-subunit  
17 ribosomal RNA genes. *Mol Biol Evol* **20**, 1881-1886,  
18 doi:10.1093/molbev/msg201 (2003).
- 19 15 Nikolaev, S. I. *et al.* The testate lobose amoebae (order Arcellinida Kent, 1880)  
20 finally find their home within Amoebozoa. *Protist* **156**, 191-202 (2005).
- 21 16 Fiore-Donno, A. M., Meyer, M., Baldauf, S. L. & Pawlowski, J. Evolution of  
22 dark-spored Myxomycetes (slime-molds): molecules versus morphology. *Mol*  
23 *Phylogenet Evol* **46**, 878-889, doi:10.1016/j.ympev.2007.12.011 (2008).
- 24 17 Berney, C. *et al.* Expansion of the 'Reticulosphere': Diversity of Novel Branching  
25 and Network-forming Amoebae Helps to Define Variosea (Amoebozoa). *Protist*  
26 **166**, 271-295, doi:10.1016/j.protis.2015.04.001 (2015).
- 27 18 Amaral Zettler, L. A. *et al.* A molecular reassessment of the leptomyxid amoebae.  
28 *Protist* **151**, 275-282 (2000).
- 29 19 Bolivar, I., Fahrni, J. F., Smirnov, A. & Pawlowski, J. SSU rRNA-based  
30 phylogenetic position of the genera *Amoeba* and *Chaos* (Lobosea,  
31 Gymnamoebia): The origin of gymnamoebae revisited. *Molecular Biology and*  
32 *Evolution* **18**, 2306-2314 (2001).
- 33 20 Smirnov, A. *et al.* Molecular phylogeny and classification of the lobose amoebae.  
34 *Protist* **156**, 129-142 (2005).
- 35 21 Cavalier-Smith, T., Chao, E. E. Y. & Oates, B. Molecular phylogeny of  
36 Amoebozoa and the evolutionary significance of the unikont *Phalansterium*.  
37 *European Journal of Protistology* **40**, 21-48 (2004).
- 38 22 Smirnov, A. V., Chao, E., Nasonova, E. S. & Cavalier-Smith, T. A revised  
39 classification of naked lobose amoebae (Amoebozoa: lobosa). *Protist* **162**, 545-  
40 570, doi:S1434-4610(11)00031-9 [pii]10.1016/j.protis.2011.04.004 (2011).
- 41 23 Adl, S. M. *et al.* Revisions to the Classification, Nomenclature, and Diversity of  
42 Eukaryotes. *J Eukaryot Microbiol* **66**, 4-119, doi:10.1111/jeu.12691 (2019).
- 43 24 Tekle, Y. I. *et al.* Phylogenomics of 'Discosea': A new molecular phylogenetic  
44 perspective on Amoebozoa with flat body forms. *Mol Phylogenet Evol* **99**, 144-  
45 154, doi:10.1016/j.ympev.2016.03.029 (2016).

- 1 25 Tekle, Y. I. & Wood, F. C. Longamoebia is not monophyletic: Phylogenomic and  
2 cytoskeleton analyses provide novel and well-resolved relationships of  
3 amoebozoan subclades. *Mol Phylogenet Evol* **114**, 249-260,  
4 doi:10.1016/j.ympev.2017.06.019 (2017).
- 5 26 Cavalier-Smith, T., Chao, E. E. & Lewis, R. 187-gene phylogeny of protozoan  
6 phylum Amoebozoa reveals a new class (Cutosea) of deep-branching,  
7 ultrastructurally unique, enveloped marine Lobosa and clarifies amoeba evolution.  
8 *Mol Phylogenet Evol* **99**, 275-296, doi:10.1016/j.ympev.2016.03.023 (2016).
- 9 27 Cavalier-Smith, T. *et al.* Multigene phylogeny resolves deep branching of  
10 Amoebozoa. *Molecular phylogenetics and evolution* **83**, 293-304,  
11 doi:10.1016/j.ympev.2014.08.011 (2015).
- 12 28 Tekle, Y. I. & Williams, J. R. Cytoskeletal architecture and its evolutionary  
13 significance in amoeboid eukaryotes and their mode of locomotion. *R Soc Open*  
14 *Sci* **3**, 160283, doi:10.1098/rsos.160283 (2016).
- 15 29 Shin, S. *et al.* Taxon sampling to address an ancient rapid radiation: a supermatrix  
16 phylogeny of early brachyceran flies (Diptera): Diptera evolution and  
17 supermatrix. *Systematic Entomology*, doi:DOI: 10.1111/syen.12275 (2017).
- 18 30 Eme, L., Sharpe, S. C., Brown, M. W. & Roger, A. J. On the age of eukaryotes:  
19 evaluating evidence from fossils and molecular clocks. *Cold Spring Harb*  
20 *Perspect Biol* **6**, doi:10.1101/cshperspect.a016139 (2014).
- 21 31 Parfrey, L. W., Lahr, D. J., Knoll, A. H. & Katz, L. A. Estimating the timing of  
22 early eukaryotic diversification with multigene molecular clocks. *Proc Natl Acad*  
23 *Sci U S A* **108**, 13624-13629, doi:10.1073/pnas.1110633108 (2011).
- 24 32 Li, Z. *et al.* Single-Copy Genes as Molecular Markers for Phylogenomic Studies  
25 in Seed Plants. *Genome Biol Evol* **9**, 1130-1147, doi:10.1093/gbe/evx070 (2017).
- 26 33 Smith, S. A., Moore, M. J., Brown, J. W. & Yang, Y. Analysis of phylogenomic  
27 datasets reveals conflict, concordance, and gene duplications with examples from  
28 animals and plants. *BMC Evol Biol* **15**, 150, doi:10.1186/s12862-015-0423-0  
29 (2015).
- 30 34 Delsuc, F., Brinkmann, H. & Philippe, H. Phylogenomics and the reconstruction  
31 of the tree of life. *Nat Rev Genet* **6**, 361-375, doi:10.1038/nrg1603 (2005).
- 32 35 Maddison, W. P. Gene trees in species trees. *Systematic Biology* **46**, 523-536  
33 (1997).
- 34 36 Michel, R. & Smirnov, A. V. The genus *Flamella* Schaeffer, 1926 (Lobosea,  
35 Gymnamoebia), with description of two new species. . *Eur J Protistol*, 400–410  
36 (1999).
- 37 37 Dykova, I., Lom, J., Dvorakova, H., Peckova, H. & Fiala, I. Didymium-like  
38 myxogastrids (class Mycetozoa) as endocommensals of sea urchins  
39 (*Sphaerechinus granularis*). *Folia Parasitol (Praha)* **54**, 1-12 (2007).
- 40 38 Fiore-Donno, A. M., Tice, A. K. & Brown, M. W. A Non-Flagellated Member of  
41 the Myxogastria and Expansion of the Echinosteliida. *J Eukaryot Microbiol* **66**,  
42 538-544, doi:10.1111/jeu.12694 (2019).
- 43 39 Kudryavtsev, A., Pawlowski, J. *Squamamoeba japonican*. g. n. sp. (Amoebozoa):  
44 a deep-sea amoeba from the Sea of Japan with a novel cell coat structure. . *Protist*  
45 **164**, 13–23 (2013).

- 1 40 Goodkov, A. V. & Seravin, L. N. Ultrastructure of the 'giant amoeba' *Pelomyxa*  
2 palustris. III. The vacuolar system; its nature, organization, dynamics and  
3 functional significance. . *Tsitologiya* **33**, 17–25 (in Russian with English  
4 summary) (1991).
- 5 41 Leadbeater BCS & Green, J. *The flagellates*. (Taylor and Francis, 2000).
- 6 42 Mitchell, D. R. The evolution of eukaryotic cilia and flagella as motile and  
7 sensory organelles. *Adv Exp Med Biol* **607**, 130-140, doi:10.1007/978-0-387-  
8 74021-8\_11 (2007).
- 9 43 Cavalier-Smith, T. Early evolution of eukaryote feeding modes, cell structural  
10 diversity, and classification of the protozoan phyla Loukozooa, Sulcozoa, and  
11 Choanozoa. *Eur J Protistol* **49**, 115-178, doi:10.1016/j.ejop.2012.06.001 (2013).
- 12 44 Cavalier-Smith, T. in *The Flagellates*. (eds S. Leadbeater & J. Green) (Taylor  
13 and Francis, 2000).
- 14 45 Stechmann, A. & Cavalier-Smith, T. Rooting the eukaryote tree by using a  
15 derived gene fusion. *Science* **297**, 89-91 (2002).
- 16 46 Cavalier-Smith, T. Only six kingdoms of life. *Proceedings of the Royal Society of*  
17 *London Series B-Biological Sciences* **271**, 1251-1262 (2004).
- 18 47 Cavalier-Smith, T. Protist phylogeny and the high-level classification of Protozoa.  
19 *European Journal of Protistology* **39**, 338-348 (2003).
- 20 48 Heiss, A. A., Walker, G. & Simpson, A. G. The flagellar apparatus of *Breviata*  
21 *anathema*, a eukaryote without a clear supergroup affinity. *Eur J Protistol* **49**,  
22 354-372, doi:10.1016/j.ejop.2013.01.001 (2013).
- 23 49 Roger, A. J. & Simpson, A. G. B. Evolution: Revisiting the Root of the Eukaryote  
24 Tree. *Current Biology* **19**, R165-R167 (2009).
- 25 50 Chistiakova, L. V., Miteva, O. A., Frolov, A. O. & Skarlato, S. O. [Comparative  
26 morphology of the subphylum Conosa Cavalier-Smith 1998]. *Tsitologiya* **55**, 778-  
27 787 (2013).
- 28 51 Derelle, R. *et al.* Bacterial proteins pinpoint a single eukaryotic root. *Proc Natl*  
29 *Acad Sci U S A* **112**, E693-699, doi:10.1073/pnas.1420657112 (2015).
- 30 52 Spiegel, F. W. in *Encyclopedia of Evolutionary Biology* (ed Richard M. Kliman)  
31 325-332 (Academic Press, 2016).
- 32 53 Fiz-Palacios, O. *et al.* Did terrestrial diversification of amoebas (amoebozoa)  
33 occur in synchrony with land plants? *PLoS One* **8**, e74374,  
34 doi:10.1371/journal.pone.0074374 (2013).
- 35 54 Schuster, J. J. & Markx, G. H. Biofilm Architecture. . *Advances in Biochemical*  
36 *Engineering/Biotechnology*, 77–96, doi:doi:10.1007/10\_2013\_248 (2013).
- 37 55 Rich, V. I. & Maier, R. M. *Aquatic Environments*. Third Edition edn, (2015).
- 38 56 Allwood, A. C., Walter, M. R., Kamber, B. S., Marshall, C. P. & Burch, I. W.  
39 Stromatolite reef from the Early Archaean era of Australia. *Nature* **441**, 714-718,  
40 doi:10.1038/nature04764 (2006).
- 41 57 Noffke, N., Christian, D., Wacey, D. & Hazen, R. M. Microbially induced  
42 sedimentary structures recording an ancient ecosystem in the ca. 3.48 billion-year-  
43 old Dresser Formation, Pilbara, Western Australia. *Astrobiology* **13**, 1103-1124,  
44 doi:10.1089/ast.2013.1030 (2013).



- 1 58 Andersson, A. *et al.* Predators and nutrient availability favor protozoa-resisting  
2 bacteria in aquatic systems. *Sci Rep* **8**, 8415, doi:10.1038/s41598-018-26422-4  
3 (2018).
- 4 59 Matz, C. & Kjelleberg, S. Off the hook--how bacteria survive protozoan grazing.  
5 *Trends Microbiol* **13**, 302-307, doi:10.1016/j.tim.2005.05.009 (2005).
- 6 60 Andersen, P. & Fenchel, T. Bacterivory by microheterotrophic flagellates in  
7 seawater samples. . *Limnol. Oceanogr.* **30**, 198–202. (1985).
- 8 61 Fenchel, T. *The Ecology of Heterotrophic Microflagellates.*, Vol. 9 (1986).
- 9 62 Purcell, E. M. Life at low Reynolds number. *American Journal of Physics* **45**  
10 (1977).
- 11 63 Fenchel, T. *Ecology of Protozoa.* (Springer-Verlag, 1987).
- 12 64 Butler, H. & Rogerson, A. Consumption rates of six species of marine benthic  
13 naked amoebae (*Gymnamoebia*) from sediments in the Clyde Sea area. . *Journal*  
14 *of the Marine Biological Association of the United Kingdom* **77**, 989-997. (1997).
- 15 65 Jackson, S. M. & Jones, E. B. G. Interactions within biofilms: the disruption of  
16 biofilm structure by protozoa. *Kieler Meeresforsch. Sonderh.* **8**, 264-268.  
17 (1991).
- 18 66 Martin, K. H., Borlee, G. I., Wheat, W. H., Jackson, M. & Borlee, B. R. Busting  
19 biofilms: free-living amoebae disrupt preformed methicillin-resistant  
20 *Staphylococcus aureus* (MRSA) and *Mycobacterium bovis* biofilms. .  
21 *Microbiology* **166**, 695. (2020).
- 22 67 Jahnke, J., Wehren, T. & Priefer, U. B. vitro studies of the impact of the naked  
23 soil amoeba *Thecamoeba similis* Greef, feeding on phototrophic soil biofilms. . *In*  
24 *Eur J Soil Biol* **43**, 14– 22 (2007).
- 25 68 Anderson, O. R. Naked amoebae in biofilms collected from a temperate  
26 freshwater pond. *J Eukaryot Microbiol* **60**, 429-431, doi:10.1111/jeu.12042  
27 (2013).
- 28 69 Rogerson, A., Anderson, O. R. & Vogel, C. re planktonic naked amoebae  
29 predominately floc associated or free in the water column? *Journal of Plankton*  
30 *Research* **25**, 1359–1365 (2003).
- 31 70 Parry, J. D. Protozoan grazing of freshwater biofilms. *Adv Appl Microbiol* **54**,  
32 167-196, doi:10.1016/S0065-2164(04)54007-8 (2004).
- 33 71 Wang, Y., Coleman-Derr, D., Chen, G. & Gu, Y. Q. OrthoVenn: a web server for  
34 genome wide comparison and annotation of orthologous clusters across multiple  
35 species. *Nucleic Acids Res* **43**, W78-84, doi:10.1093/nar/gkv487 (2015).
- 36 72 Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software  
37 version 7: improvements in performance and usability. *Mol Biol Evol* **30**, 772-  
38 780, doi:10.1093/molbev/mst010 (2013).
- 39 73 Capella-Gutierrez, S., Silla-Martinez, J. M. & Gabaldon, T. trimAl: a tool for  
40 automated alignment trimming in large-scale phylogenetic analyses.  
41 *Bioinformatics* **25**, 1972-1973, doi:10.1093/bioinformatics/btp348 (2009).
- 42 74 Nguyen, L. T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: a fast  
43 and effective stochastic algorithm for estimating maximum-likelihood  
44 phylogenies. *Mol Biol Evol* **32**, 268-274, doi:10.1093/molbev/msu300 (2015).

- 1 75 Edgar, R. C. Search and clustering orders of magnitude faster than BLAST.  
2 *Bioinformatics* **26**, 2460-2461, doi:btq461 [pii] 10.1093/bioinformatics/btq461  
3 (2010).
- 4 76 Chen, F., Mackey, A. J., Stoeckert, C. J. J. & Roos, D. S. OrthoMCL-DB:  
5 querying a comprehensive multi-species collection of ortholog groups. *Nucleic*  
6 *Acids Res.* **1**, 34, doi:doi: 10.1093/nar/gkj123. (2006).
- 7 77 Stamatakis, A., Ludwig, T. & Meier, H. RAxML-III: a fast program for maximum  
8 likelihood-based inference of large phylogenetic trees. *Bioinformatics* **21**, 456-  
9 463 (2005).
- 10 78 Zhou, X. *et al.* Quartet-Based Computations of Internode Certainty Provide  
11 Robust Measures of Phylogenetic Incongruence. *Systematic Biology* **69**, 308–324  
12 (2020).
- 13 79 Kobert, K., Salichos, L., Rokas, A., Stamatakis, A. & , C. t. I. C. a. R. M. f. P. G.  
14 T. Computing the Internode Certainty and Related Measures from Partial Gene  
15 Trees. *Molecular Biology and Evolution* **33**, 1606–1617,  
16 doi:https://doi.org/10.1093/molbev/msw040 (2016).
- 17 80 Shimodaira, H. An approximately unbiased test of phylogenetic tree selection.  
18 *Systematic Biology* **51**, 492-508 (2002).
- 19 81 Kishino, H., Miyata, T. & Hasegawa, M. Maximum likelihood inference of  
20 protein phylogeny and the origin of chloroplasts. ( ).  
21 <https://doi.org/10.1007/BF02109483>. *J Mol Evol* **31**, 151–160 (1990).

22  
23

## 24 **Acknowledgments**

25

26 This work is supported by the National Science Foundation EiR (1831958) and National  
27 Institutes of Health (1R15GM116103-02) to YIT. Additional support was gained from  
28 RSF 20-14-00195 to AS (evolutionary analysis). We would like to thank James T.  
29 Melton III, Estifanos Zerai, Ludmila Chystyakova, Sergei Karpov and Mandakini Singla  
30 for assistance in data collection, preliminary analysis and general discussions.

31

32

## 33 **Author contributions**

34

35 YIT conceived the project, led writing manuscript and helped design experiments and  
36 analysis. FW and FCW collected data, conducted analysis, and contributed to writing and  
37 editing of the manuscript. ORA and AS helped with writing, editing and organizing of the  
38 manuscript. All authors have read and approved the manuscript.

39

## 40 **Competing interests**

41

42 The authors declare that they have no competing interests.

43

44

45

46

1 **Figure captions**

2

3 **Figure 1.** Genome wide phylogeny of the Amoebozoa inferred using Maximum  
4 likelihood (ML) in IQ-TREE with LG+G4+C60+F model of evolution. The data matrix  
5 used to infer this tree consisted of 113,910 amino acid sites from the full dataset, derived  
6 from 824 genes and 113 taxa including 10 outgroup taxa. Clade supports at nodes are ML  
7 IQ-TREE 1000 ultrafast bootstrap values obtained using the same model. All branches  
8 are drawn to scale except a branch leading to Archamoebae, and *Sapocribrum*  
9 *chincoteaguense* and *Parvamoeba monoura*, that were reduced to one third and half,  
10 respectively.

11

12 **Figure 2.** A scheme illustrating the loss of kinetosomes and cilia under the different  
13 evolutionary hypotheses (A and B). Vertical hash marks on branches show loss of  
14 kinetosomes (the number lost as designated by labels on the diagram) depending on the  
15 lineage.

16

17 **Supplementary Figure caption**

18

19 **Figure S1.** Genome wide phylogeny of the Amoebozoa inferred using Maximum  
20 likelihood (ML) in IQ-TREE with LG+G4+C60+F model of evolution. The data matrix  
21 used to infer this tree consisted of 93,820 sites amino acid sites with three fast categories  
22 of sites (13%) removed from the full dataset. The data matrix consists of 824 genes and  
23 113 taxa including 10 outgroup taxa. The topology was estimated  
24 under LG+G4+C60+F+PMSF [Y1] model using a guide tree from a topology estimated  
25 using full dataset shown in Figure 1. Clade supports at nodes are ML IQ-TREE 1000  
26 ultrafast bootstrap values obtained using the same model. All branches are drawn to  
27 scale.

28

29 **Figure S2.** Maximum Likelihood tree inferred by RAxML with six fast categories of  
30 sites removed from the full dataset. The topology was estimated under  
31 PROTGAMMALG4X model. Total number of sites included after removing six fast sites  
32 categories is 70,543.

33

34 **Figure S3.** Internode certainty inferred by QuartetScores for topology in Figure 1. Values  
35 at branches are Quadripartition internode certainty (qp-ic); Lowest quartet internode  
36 certainty (lp-ic); Extended Quadripartition internode certainty (eqp-ic).

37

38 **Figure S4.** Internode certainty inferred using RAxML under PROTGAMMALG4X  
39 model for topology in Figure 1. Branch labels showed the internode certainty for a given  
40 internode with the most conflicting bipartition (left value) or all conflicting bipartitions  
41 (right value). Relative tree certainty including all conflicting bipartitions for this tree is  
42 0.978410.