

1 *Research Article submitted to Genome Biology and Evolution*

2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

A eukaryote-wide perspective on the diversity and evolution of the ARF GTPase protein family

Romana Vargová¹, Jeremy G. Wideman², Romain Derelle³, Vladimír Klimesš¹, Richard A. Kahn^{4*}, Joel B. Dacks^{5,6*}, Marek Eliáš^{1*}

¹Department of Biology and Ecology, Faculty of Science, University of Ostrava, Chittussiho 10, Ostrava, Czech Republic

²Biodesign Center for Mechanisms of Evolution, School of Life Sciences, Arizona State University, Tempe, AZ, 85287, USA

³Station d'Ecologie Théorique et Expérimentale, UMR CNRS 5321, Moulis, France

⁴Department of Biochemistry, Emory University School of Medicine, Atlanta, GA 30322

⁵Division of Infectious Disease, Department of Medicine, University of Alberta

⁶Centre for Life's Origin and Evolution, Department of Genetics, Evolution and Environment, University College of London, UK.

*Corresponding authors: rkahn@emory.edu (RAK), dacks@ualberta.ca (JBD) and marek.elias@osu.cz (ME)

35 **Abstract**

36 The evolution of eukaryotic cellular complexity is interwoven with the extensive
37 diversification of many protein families. One key family is the ARF GTPases that act in
38 eukaryote-specific processes, including membrane traffic, tubulin assembly, actin dynamics,
39 and cilia-related functions. Unfortunately, our understanding of the evolution of this family is
40 limited. Sampling an extensive set of available genome and transcriptome sequences, we
41 have assembled a dataset of over 2,000 manually curated ARF family genes from 114
42 eukaryotic species, including many deeply diverged protist lineages, and carried out
43 comprehensive molecular phylogenetic analyses. These reconstructed as many as 16 ARF
44 family members present in the last eukaryotic common ancestor (LECA), nearly doubling the
45 previously inferred ancient system complexity. Evidence for the wide occurrence and
46 ancestral origin of Arf6, Arl13 and Arl16 is presented for the first time. Moreover, Arl17,
47 Arl18 and SarB, newly described here, are absent from well-studied model organisms and as
48 a result their function(s) remain unknown. Analyses of our dataset revealed a previously
49 unsuspected diversity of membrane association modes and domain architectures within the
50 ARF family. We detail the step-wise expansion of the ARF family in the metazoan lineage,
51 including discovery of several new animal-specific family members. Delving back to its
52 earliest evolution in eukaryotes, the resolved relationship observed between the ARF family
53 paralogs sets boundaries for scenarios of vesicle coat origins during eukaryogenesis.
54 Altogether, our work fundamentally broadens the understanding of the diversity and
55 evolution of a protein family underpinning the structural and functional complexity of the
56 eukaryote cells.

57

58

59 **Key words:** ARF family, eukaryotic cell, evolution, GTPases, last eukaryotic common
60 ancestor, post-translational modifications

61

62 **Significance**

63 ARF Family GTPases are crucial regulations of a diversity of cellular compartments and
64 processes and as such the extent of this system in eukaryotes reflects both cellular complexity
65 in modern eukaryotes and its evolution. Strikingly, a comprehensive comparative genomic
66 analysis of the protein family is lacking, leaving its recent and ancient evolution poorly
67 resolved. We performed a comprehensive molecular evolutionary analysis, reconstructing a
68 highly complex ARF family complement in the Last Eukaryotic Common Ancestor,

69 including a number of paralogs never before identified as such, and we find resolved
70 relationships between the paralogs. This work has implications for cellular evolution from
71 eukaryogenesis to cellular complexity in metazoans.

72

73 **Introduction**

74 Understanding how the eukaryotic cell evolved in all its complexity is one of the greatest
75 open questions in evolutionary biology. Eukaryogenesis involved both the origin of new
76 genes and the diversification of key building blocks (Dacks et al. 2016; Eme et al. 2017).
77 Among the different building blocks, particular groups of proteins radiated early in the
78 evolution of eukaryotes and are represented by a large number of pan-eukaryotic orthologs,
79 presumably with conserved functions. One of the largest groups of proteins, acting in an
80 incredibly diverse array of cellular pathways, is the Ras superfamily of GTPases. This
81 superfamily is frequently equated with familiar and extensively studied eukaryotic “small
82 GTPases”. However, the more appropriate, i.e. evolutionary, definition conceives it as a
83 major monophyletic subgroup of the vast TRAFAC class of GTPases that also includes
84 prokaryotic representatives, larger proteins combining a Ras-related GTPase domain with
85 other functional domains, and – surprisingly to many in the field – the alpha subunits of
86 heterotrimeric G-proteins (Leipe et al. 2002). Because of its central role in so many
87 fundamental cellular functions, understanding the origin and evolution of this complex
88 superfamily of proteins is necessary for uncovering the processes by which eukaryotes
89 evolved and diversified.

90 The internal classification of the Ras superfamily is unsettled. In many overviews,
91 especially those concentrating on the eukaryotic small GTPases, the content of the
92 superfamily is pigeonholed into five major families (Ras, Rho, Rab, Ran, Arf/Sar; Colicelli
93 2004; Rojas et al. 2012), but this scheme ignores the prokaryotic superfamily members
94 (Wuichet and Sogaard-Andersen 2014), multi-domain proteins (such as the ROCO family;
95 Bosgraaf and Van Haastert 2003), and various other lineages clearly distinct from or not
96 easily classified into the well known families, such as the Gtr/Rag family (Klinger, Spang et
97 al. 2016) or RJL proteins (Elias and Archibald 2009). Understanding the diversity and the
98 evolutionary origin of the Ras superfamily in eukaryotes is a challenging task, given the
99 presence of tens to hundreds of Ras superfamily genes in each extant eukaryote genome
100 (Rojas et al. 2012). Disregarding potential (presently unknown) cases of horizontal gene
101 transfer from prokaryotic sources into particular eukaryote lineages, the wealth of Ras
102 superfamily genes in eukaryotes ultimately derives from a set of genes present in the Last

103 Eukaryote Common Ancestor (LECA). Several evolutionary analyses have attempted to
104 reconstruct LECA's complement of particular Ras superfamily subgroups and detail the
105 downstream innovation within eukaryotes. Prominent examples include analyses of the Rab
106 (Diekmann et al. 2011; Elias et al. 2012; Klöpper et al. 2012) and Ras families (van Dam et
107 al. 2011), and some isolated lineages like RJL (Elias and Archibald 2009), Miro (Vlahou et
108 al. 2011), or RABL2 (Eliáš et al. 2016). These investigations demonstrated that a large
109 number of functionally investigated paralogs were present in the LECA, emphasizing the role
110 of loss or streamlining of genomic complement in many eukaryotic lineages. They also
111 identified ancient LECA paralogs of unknown function that have been lost in lineages leading
112 to conventional model systems but which are present in diverse eukaryotic lineages of
113 ecological and medical importance. Paralogs with such an evolutionary distribution were
114 recently coined jotnarlogs (More et al. 2020). Finally, these studies also inevitably shed light
115 on the diversification of GTPases in the post-LECA expansion phase. For example, divergent
116 paralogs of unclear evolutionary relationships are found in various taxa (e.g., Pereira-Leal
117 2008), most likely resulting from rapid sequence evolution of lineage-specific paralogs linked
118 to their neofunctionalization. Additionally, the inherently small nature of the GTPases makes
119 them particularly susceptible to molecular tinkering, such as accretion of additional domains
120 or gain/loss of motifs mediating specific post-translational modifications (e.g., Záhonová et
121 al. 2018).

122 Not yet addressed in a comparable evolutionary framework is the ARF protein family.
123 This large protein family is comprised of the "true" ADP Ribosylation Factors (i.e., Arfs), as
124 well as Arf-like proteins (Arls), Arf-related protein 1 (Arfrp1), and Sar1. Clearly related are
125 the beta subunits of the signal recognition particle receptor (SR β ; Schwartz and Blobel 2003).
126 Sequence analyses have also revealed that an Arf-like ancestor, modified by insertion of a
127 novel α -helical region into its GTPase domain and high sequence divergence, gave rise to the
128 alpha subunits of heterotrimeric G-proteins (abbreviated G α ; Neuwald 2007; Anantharaman
129 et al. 2011). The distinction between Arf and Arf-like (Arl) proteins was originally made
130 based upon activity in the cholera toxin-catalyzed ADP-ribosylation of the stimulator of
131 adenylyl cyclase, G α_s , as all tested Arfs retain this functionality while the Arls did not
132 (Tamkun et al. 1991; Clark et al. 1993). However, this activity has proven of very limited
133 utility in studies of cellular functions for ARF family members as greater appreciation of both
134 the size of the family in model organisms as well as the diversity of functions became clear.
135 Thus, little if any weight should be given to whether a gene is named as an Arf, an Arl, an
136 Arfrp1, or a Sar. The ARF family is functionally heterogeneous and comprises proteins

137 involved in membrane vesicle formation (Arfs, Sar1), other aspects of vesicle traffic and
138 maintenance of membranous organelle morphology (e.g., Arl1, Arl5, or Arfrp1), microtubule
139 dynamics and mitochondrial fusion (Arl2), and cilium biogenesis and function (Arl3, Arl6,
140 Arl13) (Gillingham and Munro 2007; Donaldson and Jackson 2011; Francis et al. 2016).
141 Members of this family are critical to these diverse cellular activities and dysfunction results
142 in numerous human diseases. Family members are generally considered to be single-domain
143 small GTPases. Post-translational modifications (N-terminal myristoylation or acetylation)
144 are also often critical to the protein's localization and function.

145 An early phylogenetic study on the ARF family, limited by a lack of taxonomic
146 breadth in available genomic sequences, provided an early estimate of the ancient complexity
147 of the family in LECA and identified putative lineage-specific expansions in metazoans (Li et
148 al. 2004). The analyses showed that LECA contained at least eight ancient groups of
149 orthologs inferred from representatives being present in metazoans and at least one non-
150 opisthokont (protist or plant) eukaryote. This analysis also demonstrated that some of the
151 metazoan family members lacked close relatives in other eukaryotes, suggesting that lineage-
152 specific expansions related to metazoan multicellularity occurred. Perhaps most familiar is
153 expansion yielding the well-known and founding members of the family, Arfs 1-5. These
154 have been shown as deriving from a single ancestral gene (here referred to as Arf1 for
155 simplicity) which duplicated prior to choanoflagellates, yielding Arfs 1-3 (sometimes named
156 Class I Arfs but for convention referred to here as Arf1) and Arfs 4-5 (sometimes named
157 Class II Arfs but for convention referred to here as Arf4), with each of those diversifying into
158 five Arf paralogs around the whole genome duplications in the vertebrate lineage (Manolea et
159 al. 2010). However, since these early studies, several family members from the target species
160 (including humans) have been identified (Kahn et al. 2006) and methods of phylogenetic
161 analyses of protein sequences have advanced, including the development of the ScrollSaw
162 approach facilitating analyses of complex paralog-rich families (Elias et al. 2012). Thus, the
163 time is ripe for obtaining a much better picture of the evolution of ARF family than in the
164 previous studies.

165 To this end, we assembled, extensively curated and phylogenetically analysed a
166 dataset of ARF family sequences from a taxonomically broad selection of eukaryotic species.
167 This enabled us to revise the set of ancestral eukaryotic ARF family paralogs, which has now
168 expanded to between 14 and 16 genes. Two paralogs, described here for the first time, are not
169 represented in well studied models and point to hitherto unstudied molecular functions
170 mediated by the ARF family. We observed an unexpected diversity of domain architectures

171 challenging the dogma that ARF family proteins are only small and single-domain proteins.
172 Our analyses also unveiled a range of predicted post-translational modifications (PTMs),
173 including but not limited to well-established N-terminal myristoylation, and other molecular
174 adaptations that facilitate membrane association as a central feature of ARF family biology.
175 Finally, we identified well supported relationships between the paralogs, which have
176 implications for the inferred function of the primordial family members during
177 eukaryogenesis.

178

179 **Results and Discussion**

180 **A comprehensive dataset and phylogeny of the ARF family**

181 We first gathered all ARF family sequences (including SR β but excluding the highly
182 divergent G α proteins) from a broad diversity of eukaryotes, exploiting both publicly
183 available and privately curated genomes and transcriptomes. We did not rely solely on
184 predicted protein sequence sets but also checked the genome and transcriptome assemblies to
185 ensure maximal accuracy when it comes to statements about the absence of particular genes
186 in different taxa. All sequences were carefully validated, as described under Materials and
187 Methods, and when needed, edited (by modifications of the originally predicted gene models
188 or by changes in the assembled nucleotide sequences based on inspection of raw sequencing
189 data) to ensure maximal quality and completeness of the data. Our final dataset, provided as
190 supplementary dataset 1 (Supplementary Material online), included >2,000 manually curated
191 sequences from 114 species (supplementary table 1, Supplementary Material online). The
192 number of ARF family genes in individual species ranged from 5 in the yeast
193 *Schizosaccharomyces pombe* to 70 in the rotifer *Adineta vaga* (this high number apparently
194 reflecting the tetraploid origin of its genome; Flot et al. 2013).

195 The genes were initially annotated based on their similarity to previously
196 characterized or named ARF family genes in model organisms scored by BLAST. While this
197 procedure enabled us to recognize candidate groups of orthologs and to assign most of the
198 genes into these groups, the assignment of many sequences was uncertain or unclear and a
199 more rigorous method for establishing orthologous relationships – phylogenetic analysis of a
200 multiple sequence alignment – was required to corroborate the proposed groups of orthologs
201 and to possibly identify additional ones not readily apparent from sequence-similarity
202 comparisons. Such an analysis of the whole dataset was impractical, if not impossible, for its
203 size and the existence of divergent sequences that tend to disrupt the results of phylogenetic
204 inference. We therefore utilized of the ScrollSaw protocol previously developed to deal with

205 a similarly complex family of Rab GTPases (Elias et al. 2012) and applied by others to
206 resolve deep relationships within protein families (e.g., Vosseberg et al. 2021). This protocol
207 enables one to infer a “backbone” phylogeny of a protein family by concentrating on
208 preselected sequences likely representing slowly-evolving members of the main clades of the
209 family conserved across distantly related organismal lineages. Briefly (see Materials and
210 Methods for details), we divided the sampled species into 13 groups corresponding to major
211 eukaryotic lineages, and for each pair of groups we identified all pairs of sequences (the two
212 sequences representing the two different groups) that had mutually minimal genetic distances
213 calculated by the maximum likelihood method from a multiple sequence alignment. We then
214 gathered all the sequence pairs of all the comparisons, removed redundancies, and inferred
215 trees from the full resulting dataset (supplementary fig. 1, Supplementary Material online) or
216 after pruning sequences from selected species to further decrease the complexity of the
217 analysis (fig. 1; supplementary fig. 2, Supplementary Material online). This resulted in a
218 taxonomically rich and generally well resolved final phylogeny, which enabled us to infer
219 various aspects about the evolutionary and diversity history of the ARF family in eukaryotes.

220

221 **LECA possessed an extensive array of ARF family paralogs**

222 Dissection of the “ScrollSaw” trees indicated the existence of 13 potentially monophyletic
223 groups (Sar1 and SarB are counted as a single putative clade for the moment, see below).
224 Each group is represented by genes from all or a majority of the major eukaryote lineages, in
225 all cases spanning both putative principal clades of eukaryotes (Opimoda and Diphoda)
226 defined by the most recent hypothesis on the position of the root of the eukaryote phylogeny
227 (Derelle et al. 2015). As such these groups all are candidates for separate ARF family
228 paralogs differentiated before the radiation of extant eukaryotes and perhaps present in the
229 LECA, provided that they are monophyletic (i.e. that the root of the ARF family tree lies
230 outside of them). Our trees are inherently unrooted due to the absence of a suitable outgroup,
231 as other GTPases, including the presumably most closely related group, SR β , are too
232 divergent and their inclusion into these analyses limits the resolution of the trees. Hence, to
233 formally rule out the possibility that the root lies in any of the 13 putative clades, we
234 employed the outgroup-independent minimal ancestor deviation (MAD) method (Tria et al.
235 2017), which placed the root onto a branch separating the Arl16 group from all other groups
236 combined (fig. 1). We also note that the rooting outside any of the 13 groups implies a much
237 simpler evolutionary scenario than a root positioned into any of the groups, so hereafter we
238 treat the 13 groups as clades. Most of them have high statistical support (posterior

239 probability, SH-aLRT support, and ultrafast bootstrap values greater than or equal to 0.98,
240 98, and 98, respectively) (fig. 1). An exception is the clade denoted Arf1 and comprising
241 prototypical Arf sequences, but there is little doubt that it constitutes a coherent group of
242 orthologs. The weak signal for its monophyly may stem from a very slow evolution of Arf1
243 sequences (apparent also from very short branches in the tree) having precluded
244 accumulation of paralog-specific sequence features that would enable strong phylogenetic
245 separation from the related, more rapidly evolving (and much more strongly supported)
246 paralogs. Nevertheless, a focused analysis restricted to Arf1, Arf6, Arl1 and Arl5 allowed us
247 to use a protein alignment with more positions and recovered Arf1 as a supported
248 monophyletic clade (supplementary fig. 3, Supplementary Material online).

249 The existence of two separate clades of Arfs originated before the divergence of
250 metazoans, fungi, and plants was hypothesized previously but not convincingly demonstrated
251 (Li et al. 2004). We show that mammalian Arf6 has robustly supported orthologs in various
252 protists spanning the phylogenetic breadth of eukaryotes. The existence of a separate
253 eukaryotic Arf6 clade is further supported by comparison of intron positions in Arf genes
254 (supplementary fig. 4, Supplementary Material online). In contrast, as expected, the
255 mammalian Arf1-Arf5 proteins (class I and II Arfs) all cluster into the Arf1 clade. Our
256 analyses further demonstrate that the metazoan Arl16 has orthologs present in diverse protists
257 and thus represents a novel ancient ARF family paralog. Another previously unrecognized
258 ancient paralog, which we propose to call Arl18, was missed because it is not represented in
259 metazoans and has no characterized or named member. It is most closely related to Arl8, yet
260 the separation of Arl8 and Arl18 is apparent not only from the phylogenetic analysis (fig. 1;
261 supplementary fig. 2, Supplementary Material online) but also from their distinct exon-intron
262 structures (supplementary fig. 5, Supplementary Material online).

263 Two additional ancient eukaryotic ARF family paralogs seem to exist, although they
264 were not unambiguously supported by our phylogenetic analyses. The broader clade
265 including Sar1 proteins and their relatives has a somewhat unusual internal structure with a
266 strongly supported subclade, comprised of typical Sar1 proteins found in all taxa
267 investigated, and a more basal paraphyletic group of proteins representing different Sar1-like
268 paralogs from phylogenetically diverse protist lineages (fig. 1; supplementary fig. 2,
269 Supplementary Material online). These are not simply divergent Sar1 orthologs, as they
270 always co-occur with a *bona fide* Sar1 in each species analyzed, and multiple lines of
271 evidence suggest they constitute a separate ancient paralog of their own, which we call SarB
272 (adopting the name proposed before for a respective *Dictyostelium discoideum*

273 representative; Week et al. 2003). Specifically, some intron positions in SarB genes are
274 exclusive for this group and not shared with Sar1 (supplementary fig. 6, Supplementary
275 Material online) and the functionally important Walker B motif of SarB generally exhibits a
276 conserved tryptophan residue shared by other ARF family members and Gα proteins, as
277 opposed to a phenylalanine residue typical for Sar1 proteins (Vetter 2014; supplementary fig.
278 7, Supplementary Material online). Furthermore, a ML tree with SarB sequences constrained
279 to form a clade could not be rejected by AU test, as opposed to trees imposing topologies that
280 would correspond to the origin of SarB genes by multiple independent duplications of Sar1
281 genes proper (supplementary table 2, Supplementary Material online). Hence, it is most
282 parsimonious to interpret SarB as a *bona fide* ancient ARF family paralog different from
283 Sar1, with the phylogenetic signal for its monophyly virtually vanished over the eons. Such a
284 situation is not uncommon in phylogenetic analyses of families of short proteins with an
285 inherently limited phylogenetic signal. For instance, a similar behaviour was previously
286 observed with the highly conserved Rab1 GTPase paralog, whose undoubted monophyly was
287 also difficult to recover (Elias et al. 2012).

288 The second additional potential ancient paralog, here proposed to be called Arl17, is
289 present in various protists, certain fungi, and a single metazoan lineage, and its representative
290 contain one to three non-identical copies of a novel conserved domain C-terminal to the
291 GTPase domain (fig. 2; supplementary fig. 8, Supplementary Material online). The novel
292 ~100 residue, C-terminal domain displays no discernible homology to previously described
293 domains (even when tested by the highly sensitive HHpred searches), but occurs also in other
294 (non-Arl17) proteins from some opisthokonts and bacteria, either as a stand-alone protein
295 (e.g., EGF92317.1) or in combination with various non-GTPase domains (e.g.,
296 XP_004347279.1). Despite their unique domain architecture, no Arl17 sequences passed the
297 ScrollSaw filter, hence they are absent from the tree presented in fig. 1, and although forming
298 a clade in phylogenetic analysis, statistical support for their monophyly is lacking (fig. 2).
299 Still, the most parsimonious interpretation of our analyses is that Arl17 is an ancient ARF
300 family GTPase that was present already in LECA and had evolved from a duplication of the
301 Arf1 gene, but the tendency of the GTPase domains in Arl17 proteins to be very divergent
302 (supplementary fig. 9, Supplementary Material online) has weakened the signal for their
303 monophyly.

304 Having established the main lineages of the ARF family, we attempted to assign all
305 other genes in our full dataset (i.e. those that were excluded by the ScrollSaw protocol) into
306 them by considering sequence similarity scored by BLAST, comparison to lineage-specific

307 profile HMMs by HMMER, and by targeted phylogenetic analyses. The majority of genes in
308 our dataset could be allocated with confidence to a specific, ancient ARF family paralog,
309 enabling us to evaluate the pattern of retention of the ancient paralogs in modern eukaryotes
310 and to map the presumed gene losses to the eukaryote phylogeny (fig. 3; supplementary table
311 3, Supplementary Material online). Nevertheless, a relatively small number of genes (160 out
312 of > 2,000 sequences) remained unclassified. A majority of these likely correspond to taxon-
313 specific duplications of the standard ARF family members that have diverged substantially,
314 obscuring their actual evolutionary origin. Some cases, however, may represent excessively
315 divergent, unrecognized direct orthologs of the widespread genes. For example, several
316 unclassified genes showed potential affiliation to Arf6, yet without significant support in
317 phylogenetic analyses. These sequences all share one or more intron positions specific to
318 Arf6 (supplementary fig. 4, Supplementary Material online), supporting their annotation as
319 highly derived Arf6 genes. Future studies with a more comprehensive sampling may help
320 resolve cases such as these.

321

322 **Complex cellular repertoire inferred from the LECA complement**

323 The analyses presented above indicate that the LECA possessed at least 15 ARF family
324 genes; Arf1 and 6, Arl1, 2, 3, 5, 6, 8, 13, 16, 17, and 18, Arfrp1, Sar1, and SarB. In addition,
325 it certainly encoded SR β , excluded from our ScrollSaw analysis (hence absent from the trees
326 in fig. 1 and supplementary fig. 2, Supplementary Material online) due to its marked
327 divergence from the (core) ARF family and because SR β orthologs can be unambiguously
328 recognized by sequence similarity. Eight of these clades (Arf1, Arl1, 2, 3, 5, and 8, Arfrp1,
329 and Sar1) were previously recognized as likely ancient (Li et al. 2004) and the existence of
330 orthologs of the metazoan Arl13 in protists was also noted (e.g., Miertzschke et al. 2014),
331 although perhaps never documented by phylogenetic analyses. Our analysis thus indicates
332 that the complement of ARF family paralogs in LECA may have been twice as big as
333 previously identified, and further strengthens the idea that the LECA was a fully-fledged
334 eukaryotic cell making broad use of complex molecular machinery.

335 The cellular functions of many of the 16 ARF family GTPases in the LECA in
336 principle can be considered from what has been learned about their descendants in modern
337 eukaryotes, although our present knowledge about the function of various GTPases comes
338 from a limited number of phylogenetically biased model eukaryotes (primarily metazoans
339 and the yeast *Saccharomyces cerevisiae*, i.e. the opisthokonts) and it is not always certain to

340 what extent we can generalize from them to eukaryotes as a whole. In addition, each ARF
341 family member studied in any depth in mammalian cells has been found to act in more than
342 one pathway and typically with multiple downstream effectors (Kahn et al. 2009; Sztul et al.
343 2019), often making it difficult to assess which of these are ancient and which were acquired
344 later. Finally, we recognize that any inferences about ancient functional roles relies on an
345 assumption of functional homology across eukaryotes and an assumption of parsimonious
346 retention of pleisiomorphic traits. From a large assessment of membrane-trafficking proteins
347 that have been tested in model systems from across the eukaryotic tree, this assumption of
348 functional homology appears to be justified (Klinger et al. 2016), but does warrant being
349 explicitly named. With this caveat in mind, we summarize the key findings about the
350 different paralogs to paint a hypothetical picture of the cellular engagement of the ARF
351 family members in the LECA.

352 Most of the ARF family paralogs clearly play a role in the endomembrane dynamics.
353 As a subunit of the receptor of the signal recognition particle, SR β mediates co-translational
354 import of proteins into the ER (Schwartz and Blobel 2003). Sar1 also associates with the ER
355 and recruits subunits of the COPII coat complex to promote budding of transport vesicles
356 from the ER (Miller and Barlowe 2010). Four paralogs – Arf1, Arfrp1, Arl1 and Arl5 – are
357 physically and functionally associated with the Golgi/*trans*-Golgi network (TGN). One key
358 function of Arf1 (including the metazoan Arf1 to Arf5) is to recruit different types of vesicle
359 coats (COPI, AP-1/clathrin, AP-3) to different parts of the Golgi (Jackson and Bouvet 2014).
360 Arl1 and Arfrp1 (confusingly called Arl3p in the yeast *S. cerevisiae*) are functionally linked,
361 the latter shown to be critical for Arl1 recruitment to the *trans*-Golgi in both yeast and
362 mammalian cells (Panic et al. 2003; Setty et al. 2003; Zahn et al. 2006). Arl1 recruits several
363 effectors (e.g. golgins, arfaptins, and Arf-GEFs) to the *trans*-Golgi network (TGN) and is
364 important for endosome-to-TGN traffic (Yu and Lee 2017). The function of Arl5 is less-well
365 understood, but it may partly overlap with that of Arl1, as it also localizes to the *trans*-Golgi
366 (Houghton et al. 2012), and both the fly Arl5 and the yeast Arl1 each interact with the GARP
367 tethering complex (Panic et al. 2003; Rosa-Ferreira et al. 2015). In contrast to the Golgi
368 localizing and acting members of the ARF family, Arf6 acts predominantly at the cell surface
369 and endosomes to mediate endosome recycling, cell motility, and membrane extensions,
370 which together influence cell division, lipid/cholesterol metabolism, and changes in actin
371 dynamics (D'Souza-Schorey and Chavrier 2006; Cotton et al. 2007; Funakoshi et al. 2011;
372 Schweitzer et al. 2011). Arl8 has been implicated in controlling lysosomal motility and traffic

373 in metazoan cells (Khater et al. 2015). Its localization to the vacuolar membranes in *A.*
374 *thaliana* (Heazlewood et al. 2007) suggests that functional association of Arl8 with the
375 lysosomal/vacuolar compartment is ancestral and conserved.

376 Three paralogs, Arl3, Arl6, and Arl13 have been implicated in flagellar function
377 (Fisher et al. 2020). Arl3 has been proposed to regulate the delivery of N-myristoylated and
378 prenylated proteins to the cilium (Fansa and Wittinghofer 2016; Stephen and Ismail 2016).
379 Arl6 (also called BBS3) regulates the function of the BBSome (a protein complex involved in
380 intraflagellar transport; Mourão et al. 2014). Arl13 is involved in ciliary protein import and
381 export, purportedly mediated by its activity as a positive regulator (guanine nucleotide
382 exchange factor, GEF) for Arl3 (Gotthardt et al. 2015; Ivanova et al. 2017). Arl2 shares some
383 effectors with Arl3 and is probably involved in traffic of lipidated proteins (Van Valkenburgh
384 et al. 2001; Fansa and Wittinghofer 2016), but it has its own specific agenda, as it regulates
385 the assembly of $\alpha\beta$ -tubulin dimers (Al-Bassam 2017; Francis et al. 2017a; Francis et al.
386 2017b) and mitochondrial fusion (Newman et al. 2017).

387 Only a single study addressing the function of Arl16 has been published, reporting
388 that the mammalian Arl16 inhibits the function of the RIG-I protein, involved in the defence
389 against RNA viruses (Yang et al. 2011), but more specific functional insights are lacking.
390 Functions for of the newly discovered paralogs SarB, Arl17, and Arl18 are completely
391 unknown, as these paralogs are missing from all common model eukaryotes and thus
392 represent examples of “jotnarlogs”, proteins that are present across eukaryotes, but missing in
393 well-studied cell biological models (More et al. 2020). This adds further credence to the
394 proposal that this is a substantial evolutionary cell biological phenomenon and highlights the
395 gap in our understanding of the cell biology of the ARF family in eukaryotes. Nevertheless,
396 some clues as to the function of these proteins are provided by the phylogenetic relationship
397 to other paralogs, as relatedness within the ARF family appears to signify some level of
398 functional similarity, despite exceptions. Indeed, the aforementioned functional aspects
399 shared by the pairs Arl2-Arl3 and Arl1-Arl5 are reflected by close relationship of the
400 paralogs in the pairs (fig. 1). Likewise, the related Arf1 and Arf6 paralogs, although different
401 in terms of the intracellular localization and effectors they deploy (Jackson and Bouvet
402 2014), share the same class of GEFs and GTPase activating proteins (GAPs), though to a
403 very incompletely characterized extent (Casanova 2007; Kahn et al. 2008; Sztul et al. 2019).
404 Hence, by analogy we speculate that Arl18 may have similar functional attributes as its
405 closest paralog Arl8 (e.g. it may likewise function in the lysosomal/vacuolar sector of the
406 endomembrane system), and that SarB functions similarly to the canonical Sar1 protein in the

407 secretory pathway (Sato and Nakano 2007; Melville et al. 2020). The specific relationship of
408 Arl17 and true Arfs may be less informative concerning the function of the latter, given the
409 unique domain architecture of Arl17 proteins and the generally divergent nature of their
410 GTPase domains (compare the branch lengths of Arl17 sequences in the tree in fig. 2).

411

412 **Phylogenetic profiles of some ancestral eukaryotic ARF family paralogs illuminate**
413 **differential simplification of endomembrane system functions in eukaryote evolution**

414 A detailed scrutiny of the taxonomic distribution of some of the ancestral ARF family
415 paralogs in extant eukaryotes provides interesting insights into the variation of their roles in
416 cell functions across eukaryotes. While a hallmark of the ARF family perhaps is that
417 members are commonly found to be active in multiple, distinct pathways in the same cells
418 (Francis et al. 2016; Sztul et al. 2019), here we discuss their known or predicted
419 functionalities with respect to their best known activities, recognizing the limitations that
420 result.

421 Arfs (specifically the Arf1 paralog), Sar1, and SR β are all found in every eukaryote
422 sampled (with one exception in case of SR β , most likely due to incompleteness of the data;
423 supplementary table 3, Supplementary Material online), indicating that they belong to the
424 functional core of the eukaryotic protein toolkit. Nearly ubiquitous is Arl2, being absent only
425 from *Entamoeba histolytica*. Inspection of genomes of other *Entamoeba* species suggest that
426 Arl2 loss is not an artefact and predates the radiation of the genus. Given the role of Arl2 in
427 the assembly of tubulin dimers and in mitochondrial fusion (Francis et al. 2016), its absence
428 in *Entamoeba* may be related to a unique combination of traits of this taxon including
429 divergent tubulin sequences and a highly reduced microtubular cytoskeleton (Roy and Lohia
430 2004; Meza et al. 2006), and a simplified mitochondrion (i.e., a mitosome; Makiuchi and
431 Nozaki 2014).

432 Five of the ancestral paralogs functionally linked to the endomembrane system (based
433 on data from model eukaryotes) show various degrees of patchiness in their occurrence (fig.
434 3A; supplementary table 3, Supplementary Material online). Arl1, Arl5, and Arfrp1, all
435 associated with the Golgi apparatus, have been preserved in all main eukaryote lineages
436 sampled, but have been lost from some more terminal branches. Arl1 is missing from the
437 fission yeast (*S. pombe*), diplomonads, and some apicomplexans. Arfrp1 is absent from the
438 same set of species plus two more (the highly reduced endosymbiotic kinetoplastid
439 *Perkinsela* sp. CCAP 1560/4 and the tiny green alga *Micromonas commoda*). The similar

440 patterns of loss of these two GTPases may reflect the fact that they were shown to work in
441 the same functional cascade (see above). How Arl1 functions in the absence of Arfrp1 in
442 *Perkinsela* or *Micromonas* remains an open question but may reflect the multiplicity of
443 pathways each GTPase may influence and the potential differences in their means of
444 localization and activation. Arl5 is missing from many more eukaryotes, including even some
445 metazoans (e.g., the flatworm *Schmidtea mediterranea*). A minimum of 20 independent
446 losses of Arl5 is required to explain its distribution in our dataset (supplementary table 3,
447 Supplementary Material online), suggesting that this GTPase is a less critical component of
448 the basic infrastructure of the eukaryotic cell. In accord, disruption of the Arl5 gene in
449 *Drosophila melanogaster* does not alter the fly's viability or fertility (Rosa-Ferreira et al.
450 2015). Arl5 is closely related to Arl1 and the two GTPases may share some effectors (see
451 above). It is thus possible that Arl5 loss is facilitated by partial functional redundancy with
452 Arl1. Similar to Arl5, the distribution of Arl8 in extant eukaryotes has been shaped by
453 multiple (at least 14) independent losses, including one in the lineage leading to the main
454 eukaryotic taxon Stramenopiles (fig. 3B; supplementary table 3, Supplementary Material
455 online). Comparison of phylogenetic profiles of Arl8 and the related uncharacterized paralog
456 Arl18 reveals that the former paralog has been retained more frequently than the latter, but in
457 a few taxa (e.g., stramenopiles) Arl18 occurs in the absence of Arl8 (fig. 3A; supplementary
458 table 3, Supplementary Material online). It would be interesting to investigate whether a level
459 of functional redundancy might allow Arl18 to have taken over some of the Arl8 functions in
460 these organisms. The presence of both Arl8 and Arl18 in model systems like *Tetrahymena*
461 *thermophila* and *Trypanosoma cruzi* (supplementary table 3, Supplementary Material online)
462 provides a chance that functional dissection of these closely related paralogs is possible.

463 The patchy distribution of Arf6 is somewhat surprising, at least in part because it
464 contrasts with the near universal distribution of Arf1 paralogs. While Arf6 is perhaps most
465 commonly associated with endocytosis and plasma membrane dynamics (see above) we
466 speculate that perhaps it is its role in pericentriolar localization of specific subsets of
467 recycling endosomes that are required for midbody formation and abscission (Fielding et al.
468 2005; Wilson et al. 2005; Turn et al. 2020) that vary amongst species. The nature and
469 composition of centrioles, as well as associated components are known to vary, including
470 losses or differences in Archaeplastida and SAR (Nabais et al. 2020).

471 The unexpected discovery of the sporadically distributed, yet potentially ancestral
472 SarB paralog (figs 1 and 3A; supplementary table 3, Supplementary Material online) raises
473 an interesting possibility of a specific elaboration of the ER function in the LECA lost for

474 some reason(s) by most major eukaryotic groups. Direct functional characterization of SarB
475 in suitable model organisms is necessary before the causes behind the retention/loss pattern
476 of the gene may be understood. However, it is interesting to compare SarB with the recently
477 uncovered complexity of the ancestral set of paralogs of the COPII coat complex, including
478 the Sec24III paralog as patchily distributed as SarB (Schlacht and Dacks 2015). The
479 phylogenetic profiles of SarB and Sec24III do not overlap well (e.g., SarB is missing from
480 Chloroplastida and Sec24III is absent from diatoms), so we are not suggesting a specific
481 functional link between these two proteins. Nevertheless, the existence of both proteins
482 implies the existence of an interesting degree of variation in the COPII vesicle formation at
483 the ER in different eukaryotes.

484

485 **Arl17 provides a rare example of horizontal transfer of a Ras superfamily gene**

486 The newly recognized functionally uncharacterized Arl17 group of ARF family protein is
487 unusual not only because of its unique domain architecture (fig. 2), but also due to its very
488 patchy taxonomic distributions (fig. 3A; supplementary table 3, Supplementary Material
489 online). Based on our current sampling, Arl17 is completely missing from several major
490 eukaryotic clades (Malawimonadida, Metamonada, Discoba, Stramenopiles, Haptophyta, and
491 Rhodophyta), whereas its occurrence in the other groups is typically sporadic. Particularly
492 interesting is identification of a group of four closely related Arl17 homologs in the rotifer *A.*
493 *vaga* (Fig. 2; supplementary table 3, Supplementary Material online), which is the sole
494 representative of the densely sampled Holozoa clade possessing Arl17 (supplementary table
495 3, Supplementary Material online). Transcriptome data from *A. vaga* relatives indicate that
496 Arl17 is not restricted to a single rotifer species (data not shown), ruling out contamination in
497 the *A. vaga* genome data. Hence, the isolated occurrence of Arl17 in a rotifer lineage strongly
498 indicates gain via horizontal gene transfer (HGT) from a protist or fungal lineage, with
499 subsequent gene duplications (at least partly accounted for by tetraploidy of the *A. vaga*
500 genome, see above). Indeed, analyses of rotifer genomes revealed propensity of these
501 peculiar microscopic animals for gene gain from various sources, and three of the four *A.*
502 *vaga* Arl17 paralogs were included in the list of HGT candidates in the *A. vaga* genome (Flot
503 et al. 2013). To our knowledge, this is the first convincing case of a eukaryote-to-eukaryote
504 HGT in the whole Ras superfamily. Even though phylogenetic analysis of the Arl17 GTPase
505 domain did not shed light on the origin of rotifer's Arl17 (fig. 2), a specific relationship to
506 Arl17 proteins from *Physarum polycephalum* is suggested by a phylogeny inferred for the

507 different copies of the C-terminal novel domains (supplementary fig. 10, Supplementary
508 Material online), suggesting that rotifers acquired Arl17 from an amoebozoan.

509

510 **Expansion of the ARF family in Holozoa**

511 Given the prominent position of metazoan model systems (humans, *Mus musculus*, *D.*
512 *melanogaster*, *Caenorhabditis elegans*) in research on the ARF family, we carried out a
513 separate analysis concentrating on the family members in widely sampled representatives of
514 Metazoa and their closest protist relatives, together constituting the taxon called Holozoa.
515 Analogously to our eukaryote-scale ScrollSaw analysis described above, we compared 18
516 groups of sequences corresponding to the main holozoan lineages (phyla). This approach
517 narrowed our original holozoan dataset of nearly 550 sequences to ~320 sequences and
518 phylogenetic analysis of this reduced dataset revealed a set of strongly supported clades that
519 provided a basis for defining ARF family paralogs conserved across the main holozoan or
520 metazoan lineages (fig. 4). All ancient eukaryotic paralogs represented in this taxon, except
521 Arf1, form supported clades (note that Arl17 failed to pass the ScrollSaw step as it is present
522 only in rotifers). Furthermore, six additional groups could be identified based on this
523 analysis, namely Arf4 (class II Arf), Arl4, 10, 15, 19 and TRIM23. Most of them are named
524 according to previously annotated vertebrate genes (Gillingham and Munro 2007). An
525 exception is a novel group, here named Arl19, which is not a resolved clade, but seems to
526 represent a coherent evolutionary lineage based on additional evidence (see below). Analysis
527 of intron positions in a subgroup of ARF family genes corresponding to Arfs and their closest
528 relatives supported the delimitation of the main groups, but also suggested that several
529 sequences initially annotated as Arf1 (based on BLAST searches) may constitute a novel
530 conserved group in unicellular holozoans and several invertebrate lineages (supplementary
531 fig. 11, Supplementary Material online). Specifically, this group is characterized by three
532 unique intron positions, and a focused phylogenetic analysis supported its monophyly and
533 separation from Arf1 and other clades (supplementary fig. 12A, Supplementary Material
534 online). We thus named this novel clade Arl20.

535 Establishment of novel ARF lineages provided a basis for the assignment of
536 sequences excluded by the ScrollSaw protocol by the same approach as described for the
537 ancient eukaryotic paralogs. Moreover, inspection of the exon-intron structures facilitated
538 assignment of some of the problematic genes. For example, *Takifugu rubripes* harbours
539 several standard Arfs and one additional Arf-like paralog (TruArf4L in supplementary table
540 1, Supplementary Material online) with an almost equal similarity to the Arf1 and Arf4

541 groups. Both phylogenetic and HMMER-based analyses were inconclusive concerning the
542 origin of this gene, but the exon-intron structure of TruArf4L exhibits the pattern typical to
543 the Arf4 group (supplementary fig. 11, Supplementary Material online), supporting
544 annotation of this gene as a divergent representative of the Arf4 group. Combining such
545 different forms of evidence allowed us to annotate the majority of sequences, to establish the
546 phylogenetic distribution of the main groups, and to map their origins and losses onto the
547 holozoan phylogeny (fig. 5; supplementary table 3, Supplementary Material online).

548 Altogether we could recognise seven groups that apparently originated after the split
549 of the holozoan lineage from their relatives (Holomycota), that is in the holozoan stem itself
550 (Arf4), at a later step but still before the divergence of Metazoa and their sister group
551 choanoflagellates (Arl15, 19, 20), in the metazoan stem (Arl10), or after the divergence of the
552 deepest metazoan phyla (Arl4, TRIM23). This stepwise build-up of complexity of the ARF
553 family (fig. 5B) contrasts with a somewhat different evolutionary pattern documented for the
554 Rab family, which experienced a wave of expansion concentrated in the metazoan stem
555 lineage (Elias et al. 2012). The novel ARF family members in Holozoa apparently emerged
556 by modification of duplicated copies of specific ancient eukaryotic paralogs, although the
557 exact sources may be difficult to determine. Sequence similarity and phylogenetic analysis
558 (fig. 4) point to the Arl2/3 clade as the most likely cradle of Arl10 and 15, but the position of
559 these two paralogs is unstable in different phylogenies (e.g. supplementary fig. 12B,
560 Supplementary Material online). Evidence is more solid for the origin of Arf4, Arl4, 19, 20
561 and TRIM23, suggesting these are offshoots stemming from Arf1/6-like ancestors (fig. 4;
562 supplementary fig. 12, Supplementary Material online).

563 A common origin of Arf1 and Arf4 groups was already reported (Li et al. 2004;
564 Manolea et al. 2010), but our analysis placed this event before the divergence of
565 ichthyosporeans to the common ancestor of Holozoa (fig. 5B), which probably possessed
566 Arf1, Arf4, and Arf6 as single-copy genes. While Arf4 and Arf6 seem to duplicate only
567 sporadically, Arf1 is often present in more than one copy, suggesting a high propensity for
568 duplication; this tendency is in fact seen for eukaryote lineage in general (supplementary
569 table 3, Supplementary Material online). Phylogenetic analyses usually do not recover Arf1
570 and Arf4 as supported monophyletic clades (e.g., fig. 4), which is probably a result of their
571 high sequence similarity reflected also in partial functional overlap of Arf1 and Arf4 (Jackson
572 2014; Jackson and Bouvet 2014). However, their separation is obvious from the comparison
573 of the exon-intron structures of the respective genes (supplementary fig. 11, Supplementary
574 Material online).

575 Two more holozoan or metazoan GTPase groups are likely evolutionarily derived
576 from the ancestral Arf1 gene, yet have diverged to the point it seems inappropriate to call
577 them “Arfs”. One is Arl20, a previously unrecognized group of genes sharing three specific
578 intron positions (supplementary fig. 11, Supplementary Material online). Their relationship to
579 Arf1 cannot be conclusively inferred from our phylogenetic analysis (supplementary fig.
580 12A, Supplementary Material online), but an intron position shared with Arf1 (and Arf4) and
581 outcomes of similarity searches support this hypothesis. TRIM23 (also called ARD1) is an
582 unusual protein including not only the GTPase domain, but also a block of domains
583 characteristic for the TRIM family (RING-type E3 ubiquitin ligase, a tandem of BBbox
584 domains, and the BBC domain forming a coiled-coil) at the N-terminus. The GTPase domain
585 is highly similar to true Arfs (Vichi et al. 2005) and its specific relationship to Arf1 is
586 obvious from the virtually identical exon-intron structure (of the gene part encoding the
587 GTPase domain; supplementary fig. 11, Supplementary Material online).

588 Arl4 and the Arl19 group newly recognized here constitute a sister group to Arf6 in
589 our trees (fig. 4; supplementary fig. 12A, Supplementary Material online). While Arl4 forms
590 a highly supported monophyletic group, its placement disrupts the monophyly of Arl19,
591 perhaps due to an insufficient phylogenetic signal that would unite all Arl19 sequences in the
592 analyses. The origin of Arl4 and Arl19 from Arf6 is conceivable and there are also potential
593 functional links between Arl4 and Arf6; e.g., mammalian Arl4 proteins can recruit the Arf6
594 GEFs cytohesins to the plasma membrane (Hofmann et al. 2007) and each GTPase can
595 influence actin dynamics (Cotton et al. 2007; Li et al. 2007; Patel et al. 2011). The exon-
596 intron structure of Arl4 and Arl19 are not helpful in unveiling their origin. Only a minority of
597 Arl4 genes contain introns, the intron positions are not conserved between Arl4 genes, and do
598 not match the rest of examined Arf genes (supplementary fig. 11, Supplementary Material
599 online). This suggests that Arl4 may have originated through retroposition (Kaessmann et al.
600 2009), that is by integration of a reverse-transcribed mRNA into the genome of an early
601 metazoan, with the few non-conserved introns gained secondarily and independently in
602 different metazoan lineages. The exon-intron structure of Arl19 is rather puzzling, as several
603 genes share an intron with Arf1 (supplementary fig. 11, Supplementary Material online), but
604 the whole clade branches off close to Arf6 (fig. 4).

605 In addition to the aforementioned novel ARF family members broadly conserved
606 across Holozoa or Metazoa, various metazoan lineages exhibit still other novelties suggesting
607 further functional elaboration. Here we focus on vertebrates. First, the vertebrate ARF family
608 complement has been expanded by duplications of Arf1 and Arf4, yielding the well-known

609 two groups of paralogs (Arl1, 2, 3 versus Arf4 and 5). Together with multiple duplications of
610 Arl4, vertebrates are thus endowed with a battery of lineage-specific paralogs that are
611 generally highly similar in sequence and (presumably) function (supplementary tables 1 and
612 3, Supplementary Material online). Second, vertebrates have experienced duplication of the
613 Arl10 gene inherited from their invertebrate ancestor, giving rise to two in-paralogs that
614 diverged from each other to such an extent that they were not initially recognized as closely
615 related and which is reflected in their different names: Arl9 and Arl10 (supplementary fig.
616 12B, Supplementary Material online). Finally, vertebrates encode two divergent ARF family
617 members of a common origin, called Arl11 and Arl14, that seems to have evolved by
618 duplication and divergence from Arl4 (fig. 5; supplementary fig. 13, Supplementary Material
619 online). The functional significance of these novelties is unclear, owing to limited knowledge
620 of the function of the respective proteins in any vertebrate species including humans. It is,
621 however, important to stress that the vertebrate ARF family complement has been sculpted
622 also by gene loss, as vertebrates (in contrast to their sister group tunicates represented in this
623 study by *Ciona intestinalis*) lack Arl19 and Arl20 (fig. 5).

624

625 **The emergence of other major eukaryotic clades was accompanied by limited** 626 **evolutionary novelty in the ARF family**

627 Given the identification of multiple novel ARF family paralogs in Holozoa/Metazoa, we also
628 applied the ScrollSaw protocol to other eukaryote groups to uncover possible lineage-specific
629 innovations. Interestingly, while gene duplications specific to terminal organismal lineages
630 are common in the ARF family, only three higher-level taxa – rhodophytes, glaucophytes,
631 and Chloroplastida – seem to have evolved novel family members by gene duplication in
632 their stem lineages (fig. 3B; supplementary table 4, Supplementary Material online). The
633 genome of red algal ancestors underwent massive reductive evolution (Yoon et al. 2017),
634 which is reflected also by their highly reduced set of Rab GTPases (Petrželková and Eliáš
635 2014) as well as of ARF family proteins (fig. 3; supplementary table 3, Supplementary
636 Material online). Somewhat opposite to this trend, a novel ARF family member, here denoted
637 ArlRhodo, is shared by distantly related rhodophyte taxa and apparently emerged before the
638 radiation of the whole group. Their origin remains elusive, as the phylogenetic analysis
639 placed ArlRhodo as a separate clade of the ARF family with no specific affinities to any of
640 the ancestral clades (fig. 6A). By contrast, the glaucophyte innovation, in fact represented by
641 multiple paralogs in individual glaucophyte species, can clearly be traced as a highly
642 divergent offshoot of Arl13 (supplementary fig. 14, Supplementary Material online).

643 The only previously documented innovation of the ARF family specific for a major
644 eukaryotic group other than metazoans is the plant ArfB (Vernoud et al. 2003). It was
645 proposed to be an Arf6 ortholog (Li et al. 2004), and indeed our phylogenetic analysis places
646 ArfB as sister group to Arf6 (supplementary fig. 15, Supplementary Material online).
647 However, this topology is not statistically supported and can be an artefact resulting from the
648 apparently rapid initial evolution of the ArfB gene reflected by the long stem branch
649 subtending the ArfB subtree. Moreover, ArfB genes share one intron position with Arf1, but
650 none with Arf6 (supplementary fig. 4, Supplementary Material online). Hence, we leave the
651 origin of the ArfB group as unresolved. This notwithstanding, the timing of the ArfB
652 emergence coincides with a duplication of the ARF GEF BIG in the Chloroplastida (Pipaliya
653 et al. 2019). The duplication of the ArfB paralogs in embryophytes also coincides with the
654 duplication of GBF1 proteins in that same lineage. As both of these GEFs act on Arf1-
655 derived paralogs in metazoans at least, this lends itself to the hypothesis that ArfB is derived
656 from Arf1. It raises the further speculation that one of the BIG duplicates acts specifically on
657 ArfB in green algae and suggests that the ArfB, BIG, and GBF1 duplicates should all be
658 included in any activity assays aimed at understanding how this network functions in plant
659 cells.

660

661 **Extensive molecular tinkering in the evolution of membrane attachment mechanisms in** 662 **the ARF family**

663 It is currently understood that a large fraction of ARF family members act within
664 endomembrane traffic pathways through their actions on the surface of source membranes
665 (Gillingham and Munro 2007). This necessitates specific, and (typically) transient, membrane
666 attachment, typically relying on specific PTMs, employed by different ARF family members.
667 Our analyses illuminate the origins of the previously described means of membrane
668 association, but also finds evidence consistent with diversity in the mechanisms involved in
669 membrane association (summarized in fig. 7A-F).

670 N-terminal myristoylation (N-myristoylation) is the most common lipid modification
671 mediating the reversible membrane attachment of ARF family proteins (Kahn et al. 1988; Liu
672 et al. 2009). Two necessary prerequisites for N-myristoylation are the glycine residue at
673 the second position of the protein and specific sequence motif downstream that is recognised
674 by the myristoyl transferase catalysing the addition of the myristate moiety to the N-terminal
675 glycine (Duronio et al. 1991; Resh 1999). Once acted upon by N-myristoyl transferase, the
676 myristate group is attached through an amide bond that is permanent for the life of the

677 protein. Reversibility in membrane association is tightly linked to the activation status of the
678 ARF family protein, as the myristoylated N-terminal α -helix is accommodated in a
679 hydrophobic channel when the protein is inactive (GDP-bound) but becomes solvent exposed
680 in response to activation (GTP-binding), resulting in its propensity to bury the freed myristate
681 in a lipid bilayer (Pasqualato et al. 2002; Seidel et al. 2004; Liu et al. 2009; Liu et al. 2010).

682 Using dedicated bioinformatic tools (see Materials and Methods), we predicted this
683 post-translational modification for the majority of the proteins representing the ancestral
684 eukaryotic paralogs Arf1, Arf6, Arl1, and Arl5 (fig. 7G; supplementary tables 1 and 5,
685 Supplementary Material online), in keeping with previous experimental data from yeast and
686 mammalian proteins (Kahn et al. 1988; D'Souza-Schorey and Stahl 1995; Lee et al. 1997; Lin
687 et al. 2002). Virtually all Arf6, Arl1 and Arl5 proteins possess the conserved glycine residue
688 at the second position, and the negligible minority of those not predicted as N-myristoylation
689 targets may be false negatives. From almost 450 Arf1 genes investigated, 40 do not possess
690 the expected glycine residue and cannot be modified by myristoylation in a standard manner.
691 We note that a recent study found N-myristoyltransferase capable of acylating lysine in the
692 third position (Dian et al. 2020), though the predicting algorithms employed here did not
693 consider this possibility. Regardless, only three of the 40 Arf1 proteins without a
694 myristoylatable glycine have a lysine residue at the third position. All of them are
695 accompanied by two or more Arf1 genes that are N-myristoylated in the given organism
696 (supplementary table 1, Supplementary Material online), so they apparently represent
697 lineage-specific paralogs with a changed behaviour towards membranes. The newly
698 recognised Arl18 paralog, though not closely related to the previous four paralogs, also is
699 predicted to be ancestrally myristoylated, as all genes contain a glycine residue at the second
700 position and the majority of them are predicted as N-myristoylated (fig. 7G; supplementary
701 tables 1 and 5, Supplementary Material online). Interestingly, the Arl18 sister group Arl8
702 seems to ancestrally lack glycine at the second position (fig. 7G; supplementary table 1,
703 Supplementary Material online) and the only putatively N-myristoylated Arl8 can be found in
704 rhizarians, suggesting secondary acquisition of the myristoylation motif in this lineage. The
705 majority of Arl2 and Arl3 proteins do harbour a glycine residue at the second position, but N-
706 terminal myristoylation is predicted only for a few Arl3 proteins (fig. 7G; supplementary
707 tables 1 and 5, Supplementary Material online) and these may be false positives, considering
708 the experimental evidence for the lack of N-myristoylation in representative Arl3 proteins
709 (Sharer et al. 2002; Setty et al. 2004). The Arl6 group is clearly heterogeneous, including
710 members that certainly are not myristoylated as well as members that likely have this

711 modification. Thus, the evolutionary course leading to the distribution of N-myristoylation in
712 different ARF family members is not always clear. One possibility is an early origin of this
713 modification in an ancestor of all the clades with N-myristoylated members, followed by its
714 multiple secondary losses. However, multiple independent acquisitions is certainly a likely,
715 and mutually non-exclusive, alternative.

716 S-palmitoylation (i.e., addition of a palmitoyl moiety to one or more cysteine
717 residues) also mediates protein association with membranes, though unlike N-myristoylation
718 there are enzymes capable of reversing this acylation making it a more transient modification
719 (Zhou and Cox 2014). We again employed a suite of dedicated algorithms to predict the
720 presence of this modification in ARF family members, as described under Materials and
721 Methods. Arl15 proteins typically harbour several N-terminal cysteine residues, usually
722 predicted as S-palmitoylated (supplementary fig. 16, Supplementary Material online), and
723 approximately half of the Arl13 and Arl16 sequences analysed also contain one or more
724 putative S-palmitoylated cysteine residues in their N-terminal region (fig. 7G; supplementary
725 tables 1 and 5, Supplementary Material online). S-palmitoylation of Arl13 from *C. elegans*
726 and mammals has been confirmed experimentally and demonstrated as crucial not only for
727 the proper localization of the proteins, but also for stability and function (Cevik et al. 2010;
728 Roy et al. 2017). In a few cases, such as in the red algae-specific paralog ArlRhodo, S-
729 palmitoylation seems to accompany N-myristoylation (figs 6B and 7G; supplementary table
730 1, Supplementary Material online), similar to various other proteins, including GTPases (e.g.,
731 some G α proteins; Zhou and Cox 2014).

732 In addition to employing covalently attached saturated fatty acids, proteins also can be
733 permanently (absent proteolytic cleavage) anchored in the membrane via a transmembrane
734 domain. Of the proteins investigated here, this was previously demonstrated for SR β , a
735 protein anchored in the ER membrane via its N-terminal transmembrane region (Keenan et
736 al. 2001) that appears to be conserved in all SR β sequences investigated (fig. 7G). An N-
737 terminal transmembrane region was independently acquired by the Metazoa-specific Arl10
738 (see above) and several other ARF family members in various eukaryotes (fig. 7G;
739 supplementary table 1, Supplementary Material online). In some cases, we could confirm
740 conservation of such putative N-terminally anchored GTPases in a broader organism clade
741 beyond the species primarily targeted by our analysis, as is the case of divergent putative
742 Arf1 paralogs from *Bigelowiella natans* and other chlorarachniophytes (supplementary fig.
743 17A, Supplementary Material online) and from *Pavlova pinguis* and other haptophytes of the

744 class Pavlovophyceae (supplementary fig. 17B, Supplementary Material online). Another
745 mode of membrane attachment utilized by some ARF family members is accretion of specific
746 membrane-binding domains. This is exemplified by unusual proteins from choanoflagellates
747 and trypanosomatids that contain an N-terminal phosphoinositide-binding PH domain
748 (Lemmon 2007) connected to the ARF family GTPase domain by a long linker region (fig.
749 7E; supplementary table 1, Supplementary Material online). Finally, the eustigmatophyte
750 *Vischeria* sp. encodes a unique ARF family protein (VisArlX2 in supplementary table 1,
751 Supplementary Material online) with a long N-terminal extension lacking any detectable
752 conserved protein domain or functional motif and with a C-terminal tail ending with the
753 amino acid sequence CSIM (fig. 7F), which is reminiscent of the so-called CaaX motif (or
754 box) directing prenylation of the cysteine residue in diverse proteins (Fu and Casey 1999). A
755 similar protein, including this motif, is encoded by additional eustigmatophytes (not shown),
756 and two different prediction programs proposed the cysteine residue to be prenylated (see
757 Material and Methods for details). C-terminal prenylation is a common modification ensuring
758 membrane attachment of GTPases belonging to Rab, Ras and Rho families (Zhou and Cox
759 2014), but to our knowledge it has not been reported previously for an ARF family protein.

760 The well-studied mammalian members of the ARF family are subject to other post-
761 translational modifications (e.g., see Phosphosite Plus; <https://www.phosphosite.org/>), though
762 these either lack consensus motifs that prevent predicting their existence in other organisms
763 or have no known functional consequences, or both. One exception to this is N-terminal
764 acetylation of Arl8, which has been shown to be important for its association with lysosomal
765 membranes (Hofmann and Munro 2006). Similarly, in *S. cerevisiae* the Arfrp1 protein
766 (unfortunately named Arl3p only in this organism) is also acetylated and this is required for
767 its association with Golgi membranes (Behnia et al. 2004). Future development of
768 appropriate prediction tools, perhaps combined with dedicated biochemical investigations,
769 will be instrumental in grasping the full breath and evolutionary conservation of PTMs in the
770 ARF family.

771 In summary, the use of several different means of membrane attachment is consistent
772 with ARF family proteins acting predominantly on a membrane surface, and the diversity of
773 various membrane attachment mechanisms exhibited by this family is surprisingly extensive
774 and reminiscent of what has been described for the distantly related GTPase Rheb (Záhonová
775 et al. 2018). It is perhaps worth noting that eukaryotic organisms can vary widely in their
776 lipid composition and the same is true of different organelles in an organism, making
777 different means of membrane association likely important for this family of cell regulators

778 that most often act on membrane surfaces and can even modify the lipid composition via
779 direct activation of lipid kinases and lipases.

780

781 **Extensive diversity of multi-domain ARF family members**

782 The existence of the PH domain-containing ARF family proteins or the aforementioned
783 multi-domain TRIM23 protein (Vichi et al. 2005) counter the paradigm of ARF family
784 members being limited to single (GTPase) domain proteins with only short N- and C-terminal
785 extensions. In fact, our analyses challenge this dogma further. Although they represent a
786 minority (75 out of >2,000 sequences in our dataset), multi-domain ARF family members
787 represent a much greater number of different protein architectures involving combinations of
788 the GTPase domain of the ARF family with other functional domains than thought previously
789 (see column S in Supplementary table 1, Supplementary Material online).

790 The novel, presumably ancestral eukaryotic, Arl17 group characterized by combining
791 an Arf-related domain with varying numbers of tandemly arrayed copies of a novel
792 uncharacterized domain (fig. 2) was introduced above. Additional domain architectures are
793 found in proteins that generally seem to be lineage-specific innovations restricted to
794 particular taxa; some examples are provided in fig. 8. Similar to TRIM23, some include
795 domains linked to ubiquitination, namely the BTB domain or the F-box domain (see
796 Genschik et al. 2013), indicating recurrent recruitment of ARF family members into
797 ubiquitin-dependent regulatory circuits. Ciliates exhibit a unique protein with an ARF family
798 GTPase domain fused to a segment homologous to radial spoke protein 3 (RSP3), a
799 component of radial spokes in the axoneme (see Wirschell et al. 2008). This predicts ciliary
800 localization of this protein, and indeed, it is among the proteins detected in the ciliary
801 proteome of *T. thermophila* (Smith et al. 2005). *Entamoeba histolytica* possesses a protein
802 with a divergent C-terminal ARF family domain preceded by the VPS9 domain. The latter
803 domain is known to act as a GEF of the endosomal Rab GTPase Rab5 (Ishida et al. 2016), so
804 this protein may be part of a pathway with multiple sequentially acting GTPases similar to
805 regulatory GTPase cascades known from mammalian or yeast cells (Jones et al. 1999;
806 Mizuno-Yamasaki et al. 2012). Another unique domain combination occurs in one of the Arf
807 paralogs in the haptophyte *Emiliania huxleyi*, which is fused to the C-terminus of a block
808 including a domain of the 2OG-Fe(II) oxygenase superfamily. It is possible that the GTPase
809 domain regulates the enzyme activity of the N-terminal part of the protein. The ARF family
810 domain can combine also with other Ras superfamily GTPase domains, as demonstrated by a

811 protein from *Malawimonas californiana* with an N-terminal Rab domain and a C-terminal
812 Arf domain linked by a region containing detectable BTB and BACK domains (fig. 8A).

813 Tinkering with protein domains in ARF family proteins can be encountered in a
814 different evolutionary context than the emergence of lineage-specific paralogs. In the case of
815 Arl13, domains were acquired or lost without gene duplication, resulting in differences in
816 domain architectures between orthologous Arl13 genes. Previously characterized orthologs
817 from mammals and *Chlamydomonas reinhardtii* exhibit a poorly conserved C-terminal
818 extension that includes a region forming a coiled-coil followed by a proline-rich region (Hori
819 et al. 2008; Miertzschke et al. 2014; fig. 8B). Inspection of the large collection of Arl13
820 sequences amassed for this study revealed that this arrangement is distributed broadly across
821 the eukaryote phylogeny and likely ancestral. However, some species (represented by eleven
822 Arl13 genes out of 70 included in our dataset) depart in various way from this structure, e.g.
823 by lacking the proline-rich region or the coiled-coil. Recently, Zhang et al. (2018) identified a
824 non-canonical Arl13 gene from *Trypanosoma brucei* containing the DD_RI_PKA domain
825 (Dimerization/Docking domain of the Regulatory subunit of protein kinase A (PKA)) that is
826 essential for targeting of *T. brucei* Arl13 to the cilium. Our analysis revealed that the same
827 protein architecture is present also in *Euglena gracilis*, suggesting it is a synapomorphic
828 character for the whole Euglenozoa phylum (fig. 8B). Meanwhile, a subset of Stramenopiles
829 (oomycetes and ochrophytes) independently acquired DD_RI_PKA domain as two tandemly
830 arrayed copies (fig. 8B). DD_RI_PKA mediates interaction of PKA with A-kinase-anchoring
831 proteins (AKAPs), which regulate PKA localization in the cell (Sarma et al. 2010). Given the
832 ciliary function of Arl13 (see above), we speculate that the DD_RI_PKA domain in some
833 Arl13 proteins interacts with a cilium-localized AKAP, such as the aforementioned RSP3
834 protein (Gaillard et al. 2001; Jivan et al. 2009). In contrast, mammalian Arl13b contains the
835 simpler VxP motif in the large C-terminal domain that is required for ciliary localization
836 (Higginbotham et al. 2012; Cevik et al. 2013; Gigante et al. 2020). DD_RI_PKA domains in
837 *Phytophthora sojae* Arl13 are followed by the TUDOR domain, known for the ability to bind
838 to the methylated lysine and/or arginine residues (Botuyan and Mer 2016). The TUDOR
839 domain was independently accreted also to the C-terminus of the Arl13 from
840 *Aurantiochytrium limacinum* (fig. 8B). Another notable variant is encountered in Arl13 from
841 *B. natans* (fig. 8B) and other chlorarachniophytes (supplementary fig. 18, Supplementary
842 Material online), which exhibit a novel form of the C-terminal extension including the Ca²⁺-
843 binding EF-hand motif. Interestingly, the N-terminus of chlorarachniophyte Arl13 proteins
844 appears to be related to calcineurin B, a Ca²⁺-binding regulatory subunit of the protein

845 phosphatase calcineurin (Guerini 1997). It thus seems likely that Arl13 function is regulated
846 by Ca^{2+} in chlorarachniophytes. Exceptional is an *E. gracilis* gene (co-occurring in this
847 species with a typical Arl13 gene) that we named Arl13Triple, as it is composed of a tandem
848 triplication of a divergent Arl13-related GTPase domain (fig. 8B). The varying domain
849 architecture of Arl13 in different eukaryotes points to a substantial degree of functional
850 divergence of this key ciliary component.

851

852 **Insights into the early radiation of the ARF family**

853 In the analysis of protein family evolution, resolution between the paralogs is a tremendously
854 informative result as it allows the inference of cellular evolution of the associated organellar
855 compartments. However, such resolution has been difficult to obtain for many families. The
856 ScrollSaw methodology was a step forward in obtaining resolution for datasets with many
857 paralogs and short sequence length; e.g., Rabs and TBC proteins (Rab GAPs; Elias et al.
858 2012; Gabernet-Castello et al. 2013). Here, our application of the ScrollSaw methodology
859 also yielded a partially resolved backbone topology (fig. 1). We observed the robust
860 sisterhood of Arl8 and Arl18 and of these both to Arl16. We also observed the sisterhood of
861 Arl2 and Arl3 plus the moderately supported node uniting Arf1 with Arf6. Most notably,
862 there was a strongly resolved node grouping together Arf1, Arf6, Arl1, Arl2, Arl3 and Arl5
863 and separating them from the remainder of the paralogs.

864 This resolution provides the basis for several key inferences about the ancestral role
865 of the ARF family progenitor and some implications about the role of these proteins during
866 eukaryogenesis. Taking only the most broadly conserved biochemical and cellular features of
867 the various ARF family members, and assuming basic functional homology in orthologs to
868 their roles in LECA (Klinger et al. 2016), what is likely ancestral is a GTPase that changes
869 conformation to relocate from the cytosol to a membrane and which binds other proteins as
870 effector(s). Given the widespread role of ARF family members, this may mean a role in
871 membrane-traffic. However, with at least one resolved node separating the best-known
872 family members Arf1 and Sar1, a simple scenario of a single primordial GTPase that
873 nucleates a primordial vesicle coat-forming complex is ruled out. This suggests that the
874 proto-coatomer hypothesis (Devos et al. 2004) may well need to be modified to take a more
875 complicated scenario, including possible convergence, parallel evolution, and even merging
876 of architectures into account (Dacks and Robinson 2017; Field and Rout 2019).

877

878 **Conclusions**

879 Our comprehensive analysis of an extensive, well-curated dataset of ARF family proteins has
880 provided evolutionary insights and raised questions to be addressed by future molecular cell
881 biological exploration. The identification of 16 ancient ARF family paralogs both extends the
882 inferred complexity of LECA and sets a framework of what components can be expected to
883 be acting when delving into cellular function in diverse eukaryotes. By contrast the
884 identification of expanded complements, including novel paralogs, e.g. the metazoan Arl19
885 and Arl20, provide specific new candidates for investigation in some of the best explored and
886 heavily utilized cell biological model systems. The diversity of domain architecture
887 challenges the paradigm of this family strictly as small GTPases and begs probing of new
888 protein-protein interactions. Altogether, our work thus establishes a solid basis for future
889 more detailed investigations into the biology of ARF family proteins at a eukaryote-wide
890 scale.

891

892 **Materials and Methods**

893 **Building and curation of the ARF family dataset**

894 ARF family GTPases were searched in genome and/or transcriptome assemblies from 114
895 eukaryotic species selected such as to cover as many main eukaryotic lineages as possible
896 (sequence identifiers, source databases, and further comments are provided in supplementary
897 table 1, Supplementary Material online). The selection of taxa reflected the availability of
898 relevant data as of 2018, when the sampling was frozen to obtain a final sequence dataset for
899 all the subsequent analyses. As a result, several main eukaryote lineages, for which genome
900 or transcriptome data became available more recently (e.g. the CRuMs supergroup,
901 Telonemia, Rhodelphidia etc.), are not represented in our dataset. ARF family sequences
902 were identified using BLAST and its variants (Altschul et al. 1997). Each organism-specific
903 dataset was queried with reference members of the family and significant hits were evaluated
904 by reverse BLAST searches against an in-house extensively curated taxonomically-rich
905 database of GTPases. Query sequences being more similar to previously annotated members
906 of the ARF family (including SR β) were kept for further analysis. Existing protein sequence
907 predictions were carefully evaluated and in a many cases revised by modifying the predicted
908 exon-intron structure of the underlying gene model based on information from transcriptomic
909 data or comparison to homologous sequences. To identify genes potentially missing from
910 existing genome annotations, tblastn searches of nucleotide sequence data were carried out
911 and gene models were created anew for previously missed genes. If possible, truncated

912 sequences were completed using EST/TSA data or by iteratively recruiting raw
913 genomic/transcriptomic sequencing reads. Revised, newly predicted or extended sequences
914 are provided in supplementary table 1, Supplementary Material online. Putative pseudogenes
915 (except for the human Arf2 pseudogene sequence, which can be reconstructed) as well as
916 extremely divergent sequences with disrupted ARF family motif(s) were not included into the
917 dataset and are not listed in supplementary table 1, Supplementary Material online.

918 Each gene was initially annotated by considering results of blastp searches against our
919 comprehensive database of Ras superfamily proteins (iteratively updated by adding
920 sequences newly annotated in the course of the study). In most cases, the blastp output
921 enabled unambiguous assignment of the query sequence into one of the previously delineated
922 ortholog groups or into novel orthogroups that emerged during the study. Sequences most
923 similar to true Arfs, yet difficult to assign into the Arf1 or Arf6 groups or being visibly
924 divergent were provisionally annotated as “ArfX”. Still more divergent ARF family members
925 that did not show an apparently consistent affinity to a particular ARF family orthogroup
926 when examined by BLAST searches were provisionally annotated as “ArlX”. The annotation
927 of some of the ArfX and ArlX sequences was subsequently revised after the employment of
928 the ScrollSaw protocol described below.

929 Sequence data from the glaucophyte *Gloeochaete wittrockiana* were included despite
930 the fact that we noticed contamination of both available transcriptome assemblies
931 (MMETSP0308 and MMETSP1089; <https://www.imicrobe.us/#/projects/104>) by sequences
932 from an amoebozoan. The putative contaminants were identified by careful examination of
933 individual sequences and excluded from the dataset. Another potential contaminant (the
934 contig PCB_a545736;2 K: 25), showing a high similarity to Arl4 genes from primates, was
935 noticed in the transcriptome assembly from the breviate *Pygmaia biforma* and removed from
936 analyses.

937

938 **The ScrollSaw protocol and phylogenetic analyses**

939 A master multiple alignment was built for the identified ARF family members, excluding
940 short incomplete sequences and also all SR β sequences, as this group is noticeably different
941 from the core of the ARF family and many of its members tend to be rather divergent in their
942 sequence. Altogether, the master alignment included 1931 sequences. It was built iteratively,
943 starting with separate alignments for each group of sequences initially assigned to the same
944 (potential) orthologous group using the on-line program MAFFT (version 7), with default
945 parameters (Kato and Standley 2013). All alignments were checked by eye and further

946 edited manually using BioEdit (Hall 1999). The set of separate alignments was merged into
947 one large alignment using the on-line Merge function of MAFFT. Divergent (ArlX)
948 sequences were added to the alignment at the end. The final alignment was then manually
949 trimmed to remove poorly conserved and unreliably aligned positions. After removing
950 redundancies, the alignment comprised 1891 non-identical sequences and 148 aligned
951 positions (all falling within the GTPase domain shared across the family).

952 The alignment was subjected to an analysis essentially following the previously
953 published ScrollSaw protocol (Elias et al. 2012). The sequences were divided according to
954 the source species into 13 taxonomic groups covering the known diversity of eukaryotes:
955 Holozoa, Holomycota, Apusomonadida, Breviatea, Amoebozoa, Malawimonadida,
956 Planomonadida, Discoba, Metamonada, Archaeplastida, Cryptista, Haptista, and SAR. The
957 sizes of the groups differ substantially, since many evolutionarily important lineages were
958 represented only by a small number of species with genomic or transcriptomic data available
959 at the time when we initiated the study (in the case of Breviatea and Planomonadida by only a
960 single species). The master alignment was then subsampled by keeping only sequences from
961 each possible pair of the taxa listed above, corresponding to 78 combinations. For each of the
962 78 alignments, genetic distances between the sequences were inferred using the maximum
963 likelihood (ML) method (with the WAG+ Γ +I substitution model) implemented in Tree-
964 Puzzle 5.3 (Schmidt et al. 2002). Each resulting distance matrix was analysed using a custom
965 Python script to identify the so-called minimal-distance pairs. A minimal-distance pair
966 consists of two sequences from the two different taxonomic group compared that have
967 mutually minimal distances when distances to sequences from the other taxon are considered.
968 Minimal-distance pairs from all 78 pairwise taxon comparisons were gathered and
969 redundancies were removed, resulting in a set of 568 sequences. To further reduce the
970 complexity of the dataset we then removed all sequences that formed only one minimal-
971 distance pair in all 78 pairwise taxon comparisons combined. This step yielded the final full
972 ScrollSaw dataset comprising 354 sequences. A reduced ScrollSaw variant was prepared by
973 removing sequences from the majority of metamonads exhibiting generally divergent genes,
974 including *Monocercomonoides exilis*, *Giardia intestinalis*, *Spironucleus* spp. and
975 *Trichomonas vaginalis*.

976 The two variants of the final ScrollSaw dataset were used for inferring the ML
977 phylogenetic trees using the program IQ-TREE (Nguyen et al. 2015). The substitution model
978 (LG+I+G4) was selected by the program itself based on specific optimality criteria. Branch
979 support was assessed by the SH-aLRT test (Guindon et al. 2010) and the ultrafast bootstrap

980 approximation (Minh et al. 2013). The branch support of the reduced dataset was further
981 examined by MrBayes 3.2 (Ronquist et al. 2011) using the CIPRES Science Gateway (Miller
982 et al. 2010) with the following settings: prset aamodelpr=fixed(WAG); lset rates=gamma
983 Ngammacat=4 mcmc ngen=1000000 printfreq=10000 samplefreq=1000 nchains=4
984 burnin=80. A number of additional alignments for specific dedicated analyses, derived by
985 subsampling the master alignment or aligning the selected sequences *de novo* (using MAFFT
986 with subsequent manual editing as described above) were used for ML phylogenetic
987 inference using the same or similar approach. The alignments were in most cases trimmed
988 according to the mask applied to the master alignment. Smaller phylogenetic analysis with
989 only a subset of paralogous groups were trimmed either manually or by stand-alone version
990 of trimAl (version 1.2rev57; option automated1; Capella-Gutierrez et al. 2009) in order to
991 retrieve more positions for the ML phylogenetic analysis. The stand-alone version of IQ-
992 TREE or the IQ-TREE web server (<http://iqtree.cibiv.univie.ac.at/>; Trifinopoulos et al. 2016)
993 were used for the analyses. The substitution models were selected by the model selection
994 program implemented in the IQ-TREE. Branch support was assessed by the SH-aLRT test
995 and the ultrafast bootstrap approximation.

996 Tree topology testing was employed to test the hypothesis that SarB sequences form a
997 monophyletic group sister to the SarI group. ML trees were inferred from the reduced
998 ScrollSaw alignment with a different topological constraints (specified in supplementary table
999 2, Supplementary Material online) using IQ-TREE and the same procedure as used for
1000 computing the unconstrained tree (shown in fig. 1). The unconstrained and constrained trees,
1001 together with a sample of 1,000 trees obtained as ultrafast bootstrap replicates in the
1002 unconstrained ML search on the alignment, were then compared in IQ-TREE (-au option)
1003 with the substitution models and its parameters optimized from the original alignment (-m
1004 TEST) and using 10,000 RELL replicates. The p-values of the alternative topologies obtained
1005 with the Kishino-Hasegawa (KH), Shimodaira-Hasegawa (SH), and approximately unbiased
1006 (AU) tests were considered.

1007

1008 **Annotation of sequences**

1009 The full and reduced ScrollSaw datasets were used as a basis for annotation of the rest of the
1010 sequences. The identity of individual sequences or their groups was tested by adding them to
1011 the reduced ScrollSaw dataset and inferring a ML tree with IQ-TREE. The scrutinized
1012 sequences were assigned to a particular ancestral eukaryotic paralog and annotated
1013 accordingly if they clustered together with reference representatives of the given paralog

1014 group and the relationship was supported by SH-aLRT and ultrafast bootstrap values of ≥ 80
1015 and 95, respectively. Not all genes could be annotated by this approach, hence the HMMER
1016 package (stand alone version 3.0; hmmer.org) was employed as an alternative. The aligned
1017 full ScrollSaw dataset was divided into 14 separate alignments, each representing one
1018 ancestral paralog (Arf1, 6, Arl1, 2, 3, 5, 6, 8, 13, 16, 18, Arfrp1, Sar1, and SarB). A profile
1019 HMM was constructed for each alignment using hmmbuild and a database of profile HMMs
1020 was created using hmmcompress. The unannotated sequences were then used as queries in
1021 hmmscan searches against the database and the “best 1 domain” score difference between the
1022 first and the second best hits was determined. If this difference was equal to or higher than
1023 20, the sequence was annotated according to the best hit. Sequences annotated based on the
1024 phylogenetic analyses or hmmscan searches are marked by asterisk (*) in the column
1025 “Conclusively annotated” in the supplementary table 1, Supplementary Material online.
1026 Proteins representing SR β and Arl17, which were not represented by reference sequences in
1027 the ScrollSaw dataset, were unequivocally identified owing to the distinct characters of these
1028 sequence groups, which makes them easy to recognise by BLAST-based similarity searches
1029 (SR β) or by considering the presence of the novel conserved C-terminal domain (Arl17; see
1030 the main text). All SR β and Arl17 proteins are therefore also considered as conclusively
1031 annotated. A single truncated sequence (Arl17b gene from *Chromera velia*) lacked the C-
1032 terminal extension with the characteristic C-terminal domain, but was assigned to the Arl17
1033 group based on its close sequence similarity to undisputed Arl17 sequences. A combination
1034 of BLAST searches, ML phylogenetic analyses and comparison of exon-intron structure was
1035 used to obtain the most likely annotation of the sequences that could not be conclusively
1036 annotated by the aforementioned approaches. Several sequences were annotated as Arf1/6, as
1037 they showed affinity to the Arf1/6 clade, but it was impossible to decide whether they
1038 originated from ancestral Arf1 or Arf6 paralogs. Only 160 out of more than 2000 ARF family
1039 sequences analysed could not be annotated with any confidence, so they remained unassigned
1040 to any ancestral paralog (supplementary tables 1 and 3, Supplementary Material online).

1041

1042 **Taxon-specific ScrollSaw analyses and annotation of lineage-specific paralogs**

1043 To detect lineage-specific paralogs, we applied the ScrollSaw protocol separately to sets of
1044 sequences from the following main eukaryote taxa: Chloroplastida, Rhodophyta,
1045 Glaucophyta, Cryptista, Haptista, SAR, Discoba, Metamonada, Amoebozoa, Holomycota and
1046 Holozoa. Lineages represented by only one or two species (Apusomonadida, Breviatea,
1047 Planomonadida, and Malawimonadida) were not included. The ScrollSaw protocol and

1048 phylogenetic analyses of the resulted datasets were performed generally as described above
1049 for the whole dataset. For each main eukaryote taxon analysed, species representing it were
1050 assigned to predefined monophyletic subgroups specified in supplementary table 4,
1051 Supplementary Material online. Sequences from these species were extracted from the
1052 trimmed master alignment of the ARF family protein (except for sequences from Holozoa,
1053 which were aligned *de novo* using MAFFT and then trimmed according to the mask used for
1054 the whole dataset), the ScrollSaw protocol was applied to identify minimal-distance pairs,
1055 and ML phylogenetic trees were calculated on the filtered sequences. In contrast to the pan-
1056 eukaryotic ScrollSaw analysis, sequences that formed only one minimal-distance pair were
1057 not omitted (except for the analysis of the Holozoa dataset, where the criterion of the
1058 sequence belonging to at least two minimal-distance pairs was kept). The ML trees were
1059 inspected to identify robustly supported clades that would define conserved paralogs
1060 ancestral for the focal eukaryotic taxon but different from the previously defined ancestral
1061 eukaryote paralogs. In the case of Holozoa, paralogs specific for individual subgroups were
1062 considered, too. Further representatives of these paralogs (i.e., specific orthologs of the
1063 constituent sequences identified in the ScrollSaw trees) were then identified among the
1064 sequences that did not pass the ScrollSaw step by a combination of BLAST searches,
1065 phylogenetic analyses and (in case of Holozoa) HMMER-based comparisons. Candidates for
1066 ancestral taxon-specific paralogs were detected only in Chloroplastida, Rhodophyta, and
1067 Glaucophyta, as described in detail in the main text.

1068

1069 **Prediction of transmembrane regions and post-translation modifications**

1070 The presence of transmembrane (TM) regions in ARF family proteins was examined using
1071 the online TMHMM Server v. 2.0 (<http://www.cbs.dtu.dk/services/TMHMM/>). In the case of
1072 sequences with suspicious TM absence or presence (i.e., when the result was untypical for the
1073 respective ARF family subgroup), the on-line tool TMPred
1074 (https://www.ch.embnet.org/software/TMPRED_form.html) was additionally employed. The
1075 predictions are listed in supplementary table 1, Supplementary Material online. N-terminal
1076 myristoylation of sequences with the glycine residue at the second position was evaluated
1077 using the on-line ExpASy Myristoylator tool (<http://web.expasy.org/myristoylator/>; Bologna
1078 et al. 2004), NMT - The MYR Predictor
1079 (<http://mendel.imp.ac.at/myristate/SUPLpredictor.htm>; Maurer-Stroh et al. 2002), and the
1080 stand-alone version of GPS-Lipid (v1.0, <http://lipid.biocuckoo.org/index.php>; Xie et al.
1081 2016). Only those proteins predicted as N-terminally myristoylated by at least two tools were

1082 considered as significant candidates. Setting of all tools was default except for NMT - The
1083 MYR Predictor where only N-terminal glycine residues were considered, and fungal
1084 sequences were predicted with the “Fungi specific” option. In case of GPS-Lipid, the
1085 threshold was set to “low”. Possible S-palmitoylation was predicted using SeqPalm
1086 (<http://lishuyan.lzu.edu.cn/seqpalm/>; Li et al. 2015), the stand-alone version of CKSAAP-
1087 Palm programme (<http://doc.aporc.org/wiki/CKSAAP-Palm>; Wang et al. 2009), PalmPred
1088 (<http://proteininformatics.org/mkumar/palmpred/index.html>; Kumari et al. 2014), stand-alone
1089 version of GPS-Lipid, and WAP-Palm (<http://bioinfo.ncu.edu.cn/WAP-Palm.aspx>; Shi et al.
1090 2013). Only those sites predicted as S-palmitoylated by at least three tools were considered as
1091 significant candidates. Setting of all tools was default except for GPS-Lipid with the
1092 threshold set to “high”. Complete results from all tools are showed in supplementary table 5,
1093 Supplementary Material online, consensual results are included in supplementary table 1,
1094 Supplementary Material online. Possible prenylation was assessed only for VisArlX2 from
1095 the alga *Vischeria* sp., as it is the only protein from our dataset with a typical C-terminal
1096 prenylation motif. The online programs iPreny-PseAAC ([http://app.aporc.org/iPreny-
1097 PseAAC/index.html](http://app.aporc.org/iPreny-PseAAC/index.html); Xu et al. 2017) and GPS-Lipid were used with default settings; both
1098 tools predicted VisArlX2 as a prenylated protein.

1099

1100 **Other sequence analyses**

1101 Intron positions were investigated in four groups of ARF family genes (Sar1/SarB;
1102 Arl8/Arl18; Arfs and the GTPase domain of Arl17; Arfs and selected Arf-like in Holozoa) as
1103 a means to illuminate the origin and relationships of these genes. The positions of introns
1104 (including their phases) were mapped onto a multiple alignment of respective protein
1105 sequences using a custom Java script. The multiple sequence alignments were constructed *de*
1106 *novo* using MAFFT, inspected visually and adjusted manually whenever necessary
1107 (Sar1/SarB, Arl8/Arl18, Arf, and Arf-like in holozoans). For the analysis of Arf and Arl17
1108 genes, the respective protein sequences were extracted from the master alignment. Sequences
1109 with no introns in the coding sequence or represented only by transcriptomic data were
1110 omitted. A manually curated dataset of gene exon-intron structures was used as the input for
1111 the intron positions mapping. For presentation purposes, regions corresponding to
1112 unconserved N- and C- termini of the sequences were trimmed and long sequence-specific
1113 insertions were collapsed. To highlight the pattern of protein sequence conservation,
1114 CHROMA (ver. 1.0 Goodstadt and Ponting 2001) was used for processing some of the
1115 multiple sequence alignments presented. Sequence logos of the Walker B motif were

1116 obtained using the on-line tool WebLogo 3 (<http://weblogo.threeplusone.com/create.cgi>;
1117 Crooks et al. 2004) from the multiple sequence alignment of the respective sequences after
1118 removing sequence-specific insertions present in a few sequences.

1119 Conserved protein domains and other structural features in ARF family proteins were
1120 identified using searches of Pfam (<http://pfam.xfam.org/>; Finn et al. 2016), the Conserved
1121 Domains database (<https://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi>; Marchler-Bauer
1122 et al. 2017), and the SMART database (<http://smart.embl-heidelberg.de/>; Letunic and Bork
1123 2018). Domain predictions provided by the three tools were compared and spurious results
1124 (low-significance with only a single tool) were ignored. The identity of unusual N-terminal
1125 extensions present in some Arl13 proteins were evaluated using HHpred
1126 (<https://toolkit.tuebingen.mpg.de/>; Söding et al. 2005). Multiple sequence alignments of the
1127 different forms on the N-terminal extensions conserved within different taxa were used as
1128 queries in the HHpred searches. In case of the N-terminal extension conserved in Arl13
1129 proteins from Euglenozoa, the sampling was expanded beyond the focal set of taxa (including
1130 only three euglenozoans) by adding to the alignment several additional euglenozoan Arl13
1131 sequences to improve the representativeness of the alignment. Similarly, additional
1132 chlorarachniophyte Arl13 sequences were identified and aligned with the sole representative
1133 in the focal dataset (that from *B. natans*), and additional stramenopile (oomycete and
1134 ochrophyte) Arl13 sequences with the same conserved N-terminal extension as the
1135 stramenopile sequences in the focal set were included to increase the sensitivity of the
1136 analysis. Some TRIM23 sequences were predicted by the standard tools to contain only one
1137 BBOX domain rather than the two common in most members of this group, but inspection of
1138 a multiple sequence alignment revealed high similarity of all sequences in the respective
1139 region, suggesting that all TRIM23 sequences likely conform to the same domain architecture
1140 with two BBOX domains.

1141

1142 **Supplementary Material**

1143 Supplementary data are available at Genome Biology and Evolution online.

1144

1145 **Acknowledgements**

1146 We thank Eunsoo Kim (American Museum of Natural History, New York) for an access to a
1147 genome assembly of *Goniomonas avonlea* prior to publication and Vladimír Hampl (Charles
1148 University in Prague) for his permission to use sequences from an unpublished genome
1149 assembly of *Paratrimastix pyriformis*. This work was supported by Czech Science

1150 Foundation grant 20-27648S; ERD Funds, project OPVVV
1151 CZ.02.1.01/0.0/0.0/16_019/0000759 (Centre for research of pathogenicity and virulence of
1152 parasites); and the infrastructure grant CZ.1.05/2.1.00/19.0388 („Přístroje IET“). Work in the
1153 Kahn lab is supported by a grant from the National Institutes of Health (NIH
1154 R35GM122568). Work in the Dacks Lab is supported by grants from the Natural Sciences
1155 and Engineering Research Council of Canada (RES0021028, RES0043758, RES0046091).
1156 JBD is the Canada Research Chair (Tier II) in Evolutionary Cell Biology.

1157

1158 **Data availability**

1159 ARF family gene sequences extracted from unpublished genome assemblies of *Gefionella*
1160 *okellyi*, *Planomonas micra*, and *Paratrimastix pyriformis*, and from our unpublished
1161 transcriptome assembly of *Vicheria* sp. CAUP Q 202, were deposited at GenBank with
1162 accession numbers #####-#####. A complete set of manually curated eukaryotic ARF family
1163 protein sequences is available in supplementary dataset 1, Supplementary Material online.

1164

1165 **Competing interests**

1166 No competing interests declared.

1167

1168 **Literature Cited**

- 1169 Al-Bassam, J. (2017). Revisiting the tubulin cofactors and Arl2 in the regulation of soluble
1170 $\alpha\beta$ -tubulin pools and their effect on microtubule dynamics. *Molecular Biology of the*
1171 *Cell*, 28(3), 359–363. <https://doi.org/10.1091/mbc.E15-10-0694>
1172 Al-Bassam J. 2017. Revisiting the tubulin cofactors and Arl2 in the regulation of soluble $\alpha\beta$ -
1173 tubulin pools and their effect on microtubule dynamics. *Mol Biol Cell*. 28:359–363. doi:
1174 10.1091/mbc.E15-10-0694.
1175 Altschul SF et al. 1997. Gapped BLAST and PSI-BLAST: a new generation of protein
1176 database search programs. *Nucleic Acids Res*. 25:3389–3402. doi:
1177 10.1093/nar/25.17.3389.
1178 Anantharaman V, Abhiman S, de Souza RF, Aravind L. 2011. Comparative genomics
1179 uncovers novel structural and functional features of the heterotrimeric GTPase signaling
1180 system. *Gene*. 475:63–78. doi: 10.1016/j.gene.2010.12.001.
1181 Barlow LD, Nývltová E, Aguilar M, Tachezy J, Dacks JB. 2018. A sophisticated,
1182 differentiated Golgi in the ancestor of eukaryotes. *BMC Biol*. 16:27. doi:
1183 10.1186/s12915-018-0492-9.
1184 Behnia R, Panic B, Whyte JRC, Munro S. 2004. Targeting of the Arf-like GTPase Arl3p to
1185 the Golgi requires N-terminal acetylation and the membrane protein Sys1p. *Nat Cell Biol*.
1186 6:405–413. doi: 10.1038/ncb1120.
1187 Bologna G, Yvon C, Duvaud S, Veuthey A-L. 2004. N-Terminal myristoylation predictions
1188 by ensembles of neural networks. *Proteomics*. 4:1626–1632. doi:
1189 10.1002/pmic.200300783.

- 1190 Bosgraaf L, Van Haastert PJM. 2003. Roc, a Ras/GTPase domain in complex proteins.
1191 Biochim Biophys Acta. 1643:5–10. doi: 10.1016/j.bbamcr.2003.08.008.
- 1192 Botuyan MV, Mer G. 2016. Chapter 8 - Tudor Domains as Methyl-Lysine and Methyl-
1193 Arginine Readers. In: Chromatin Signaling and Diseases. Binda, O & Fernandez-Zapico,
1194 ME, editors. Academic Press: Boston pp. 149–165. doi: 10.1016/B978-0-12-802389-
1195 1.00008-3.
- 1196 Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. 2009. trimAl: a tool for automated
1197 alignment trimming in large-scale phylogenetic analyses. Bioinformatics. 25:1972–1973.
1198 doi: 10.1093/bioinformatics/btp348.
- 1199 Casanova JE. 2007. Regulation of Arf activation: the Sec7 family of guanine nucleotide
1200 exchange factors. Traffic. 8:1476–1485. doi: 10.1111/j.1600-0854.2007.00634.x.
- 1201 Cevik S et al. 2013. Active transport and diffusion barriers restrict Joubert Syndrome-
1202 associated ARL13B/ARL-13 to an Inv-like ciliary membrane subdomain. PLoS Genet.
1203 9:e1003977. doi: 10.1371/journal.pgen.1003977.
- 1204 Cevik S et al. 2010. Joubert syndrome Arl13b functions at ciliary membranes and stabilizes
1205 protein transport in *Caenorhabditis elegans*. J Cell Biol. 188:953–969. doi:
1206 10.1083/jcb.200908133.
- 1207 Clark J et al. 1993. Selective amplification of additional members of the ADP-ribosylation
1208 factor (ARF) family: cloning of additional human and *Drosophila* ARF-like genes. Proc
1209 Natl Acad Sci U S A. 90:8952–8956. doi: 10.1073/pnas.90.19.8952.
- 1210 Colicelli J. 2004. Human RAS superfamily proteins and related GTPases. Sci STKE.
1211 2004:RE13. doi: 10.1126/stke.2502004re13.
- 1212 Cotton M et al. 2007. Endogenous ARF6 interacts with Rac1 upon angiotensin II stimulation
1213 to regulate membrane ruffling and cell migration. Mol Biol Cell. 18:501–511. doi:
1214 10.1091/mbc.e06-06-0567.
- 1215 Crooks GE, Hon G, Chandonia J-M, Brenner SE. 2004. WebLogo: A Sequence Logo
1216 Generator. Genome Res. 14:1188–1190. doi: 10.1101/gr.849004.
- 1217 Dacks JB et al. 2016. The changing view of eukaryogenesis - fossils, cells, lineages and how
1218 they all come together. J. Cell. Sci. 129:3695–3703. doi: 10.1242/jcs.178566.
- 1219 Dacks JB, Robinson MS. 2017. Outerwear through the ages: evolutionary cell biology of
1220 vesicle coats. Curr Opin Cell Biol. 47:108–116. doi: 10.1016/j.ceb.2017.04.001.
- 1221 van Dam TJP, Bos JL, Snel B. 2011. Evolution of the Ras-like small GTPases and their
1222 regulators. Small GTPases. 2:4–16. doi: 10.4161/sgtp.2.1.15113.
- 1223 Derelle R et al. 2015. Bacterial proteins pinpoint a single eukaryotic root. Proc Natl Acad Sci
1224 U S A. 112:E693–E699. doi: 10.1073/pnas.1420657112.
- 1225 Devos D et al. 2004. Components of coated vesicles and nuclear pore complexes share a
1226 common molecular architecture. PLoS Biol. 2:e380. doi: 10.1371/journal.pbio.0020380.
- 1227 Dian C et al. 2020. High-resolution snapshots of human N-myristoyltransferase in action
1228 illuminate a mechanism promoting N-terminal Lys and Gly myristoylation. Nat Commun.
1229 11:1132. doi: 10.1038/s41467-020-14847-3.
- 1230 Diekmann Y et al. 2011. Thousands of Rab GTPases for the Cell Biologist. PLOS Comput
1231 Biol. 7:e1002217. doi: 10.1371/journal.pcbi.1002217.
- 1232 Donaldson JG, Jackson CL. 2011. Arf Family G Proteins and their regulators: roles in
1233 membrane transport, development and disease. Nat Rev Mol Cell Biol. 12:362–375. doi:
1234 10.1038/nrm3117.
- 1235 D'Souza-Schorey C, Chavrier P. 2006. ARF proteins: roles in membrane traffic and beyond.
1236 Nat Rev Mol Cell Biol. 7:347–358. doi: 10.1038/nrm1910.
- 1237 D'Souza-Schorey C, Stahl PD. 1995. Myristoylation is required for the intracellular
1238 localization and endocytic function of ARF6. Exp Cell Res. 221:153–159. doi:
1239 10.1006/excr.1995.1362.

- 1240 Duronio RJ, Rudnick DA, Adams SP, Towler DA, Gordon JI. 1991. Analyzing the substrate
1241 specificity of *Saccharomyces cerevisiae* myristoyl-CoA:protein N-myristoyltransferase
1242 by co-expressing it with mammalian G protein α subunits in *Escherichia coli*. J Biol
1243 Chem. 266:10498–10504. doi: 10.1016/S0021-9258(18)99252-5.
- 1244 Elias M, Archibald JM. 2009. The RNL family of small GTPases is an ancient eukaryotic
1245 invention probably functionally associated with the flagellar apparatus. Gene. 442:63–72.
1246 doi: 10.1016/j.gene.2009.04.011.
- 1247 Elias M, Brighthouse A, Gabernet-Castello C, Field MC, Dacks JB. 2012. Sculpting the
1248 endomembrane system in deep time: high resolution phylogenetics of Rab GTPases. J
1249 Cell Sci. 125:2500–2508. doi: 10.1242/jcs.101378.
- 1250 Eliáš M, Klimeš V, Derelle R, Petrželková R, Tachezy J. 2016. A paneukaryotic genomic
1251 analysis of the small GTPase RABL2 underscores the significance of recurrent gene loss
1252 in eukaryote evolution. Biol Direct. 11:5. doi: 10.1186/s13062-016-0107-8.
- 1253 Eme L, Spang A, Lombard J, Stairs CW, Ettema TJG. 2017. Archaea and the origin of
1254 eukaryotes. Nat Rev Microbiol. 15:711–723. doi: 10.1038/nrmicro.2017.133.
- 1255 Fansa EK, Wittinghofer A. 2016. Sorting of lipidated cargo by the Arl2/Arl3 system. Small
1256 GTPases. 7:222–230. doi: 10.1080/21541248.2016.1224454.
- 1257 Field MC, Rout MP. 2019. Pore timing: the evolutionary origins of the nucleus and nuclear
1258 pore complex. F1000Res. 8. doi: 10.12688/f1000research.16402.1.
- 1259 Fielding AB et al. 2005. Rab11-FIP3 and FIP4 interact with Arf6 and the exocyst to control
1260 membrane traffic in cytokinesis. EMBO J. 24:3389–3399. doi:
1261 10.1038/sj.emboj.7600803.
- 1262 Finn RD et al. 2016. The Pfam protein families database: towards a more sustainable future.
1263 Nucleic Acids Res. 44:D279-285. doi: 10.1093/nar/gkv1344.
- 1264 Fisher S, Kuna D, Caspary T, Kahn RA, Sztul E. 2020. ARF family GTPases with links to
1265 cilia. Am J Physiol Cell Physiol. 319:C404–C418. doi: 10.1152/ajpcell.00188.2020.
- 1266 Flot J-F et al. 2013. Genomic evidence for ameiotic evolution in the bdelloid rotifer *Adineta*
1267 *vaga*. Nature. 500:453–457. doi: 10.1038/nature12326.
- 1268 Francis JW, Goswami D, et al. 2017. Nucleotide Binding to ARL2 in the TBCD·ARL2· β -
1269 Tubulin Complex Drives Conformational Changes in β -Tubulin. J Mol Biol. 429:3696–
1270 3716. doi: 10.1016/j.jmb.2017.09.016.
- 1271 Francis JW, Newman LE, Cunningham LA, Kahn RA. 2017. A Trimer Consisting of the
1272 Tubulin-specific Chaperone D (TBCD), Regulatory GTPase ARL2, and β -Tubulin Is
1273 Required for Maintaining the Microtubule Network. J Biol Chem. 292:4336–4349. doi:
1274 10.1074/jbc.M116.770909.
- 1275 Francis JW, Turn RE, Newman LE, Schiavon C, Kahn RA. 2016. Higher order signaling:
1276 ARL2 as regulator of both mitochondrial fusion and microtubule dynamics allows
1277 integration of 2 essential cell functions. Small GTPases. 7:188–196. doi:
1278 10.1080/21541248.2016.1211069.
- 1279 Fu HW, Casey PJ. 1999. Enzymology and biology of CaaX protein prenylation. Recent Prog
1280 Horm Res. 54:315–342; discussion 342-343.
- 1281 Funakoshi Y, Hasegawa H, Kanaho Y. 2011. Regulation of PIP5K activity by Arf6 and its
1282 physiological significance. J Cell Physiol. 226:888–895. doi: 10.1002/jcp.22482.
- 1283 Gabernet-Castello C, O'Reilly AJ, Dacks JB, Field MC. 2013. Evolution of Tre-
1284 2/Bub2/Cdc16 (TBC) Rab GTPase-activating proteins. Mol Biol Cell. 24:1574–1583.
1285 doi: 10.1091/mbc.E12-07-0557.
- 1286 Gaillard AR, Diener DR, Rosenbaum JL, Sale WS. 2001. Flagellar radial spoke protein 3 is
1287 an A-kinase anchoring protein (AKAP). J Cell Biol. 153:443–448. doi:
1288 10.1083/jcb.153.2.443.

- 1289 Genschik P, Sumara I, Lechner E. 2013. The emerging family of CULLIN3-RING ubiquitin
1290 ligases (CRL3s): cellular functions and disease implications. *EMBO J.* 32:2307–2320.
1291 doi: 10.1038/emboj.2013.173.
- 1292 Gigante ED, Taylor MR, Ivanova AA, Kahn RA, Caspary T. 2020. ARL13B regulates Sonic
1293 hedgehog signaling from outside primary cilia. *Elife.* 9. doi: 10.7554/eLife.50434.
- 1294 Gillingham AK, Munro S. 2007. The small G proteins of the Arf family and their regulators.
1295 *Annu. Rev. Cell Dev. Biol.* 23:579–611. doi: 10.1146/annurev.cellbio.23.090506.123209.
- 1296 Goodstadt L, Ponting CP. 2001. CHROMA: consensus-based colouring of multiple
1297 alignments for publication. *Bioinformatics.* 17:845–846. doi:
1298 10.1093/bioinformatics/17.9.845.
- 1299 Gotthardt K et al. 2015. A G-protein activation cascade from Arl13B to Arl3 and
1300 implications for ciliary targeting of lipidated proteins. *Elife.* 4. doi: 10.7554/eLife.11859.
- 1301 Guerini D. 1997. Calcineurin: not just a simple protein phosphatase. *Biochem Biophys Res*
1302 *Commun.* 235:271–275. doi: 10.1006/bbrc.1997.6802.
- 1303 Guindon S et al. 2010. New Algorithms and Methods to Estimate Maximum-Likelihood
1304 Phylogenies: Assessing the Performance of PhyML 3.0. *Syst Biol.* 59:307–321. doi:
1305 10.1093/sysbio/syq010.
- 1306 Hall TA. 1999. BioEdit: a user-friendly biological sequence alignment editor and analysis
1307 program for Windows 95/98/NT. *Nucl. Acids. Symp. Ser.* 41:95–98.
- 1308 Heazlewood JL, Verboom RE, Tonti-Filippini J, Small I, Millar AH. 2007. SUBA: the
1309 *Arabidopsis* Subcellular Database. *Nucleic Acids Res.* 35:D213–D218. doi:
1310 10.1093/nar/gkl863.
- 1311 Higginbotham H et al. 2012. Arl13b in primary cilia regulates the migration and placement of
1312 interneurons in the developing cerebral cortex. *Dev Cell.* 23:925–938. doi:
1313 10.1016/j.devcel.2012.09.019.
- 1314 Hofmann I, Munro S. 2006. An N-terminally acetylated Arf-like GTPase is localised to
1315 lysosomes and affects their motility. *J Cell Sci.* 119:1494–1503. doi: 10.1242/jcs.02958.
- 1316 Hofmann I, Thompson A, Sanderson CM, Munro S. 2007. The Arl4 family of small G
1317 proteins can recruit the cytohesin Arf6 exchange factors to the plasma membrane. *Curr*
1318 *Biol.* 17:711–716. doi: 10.1016/j.cub.2007.03.007.
- 1319 Hori Y, Kobayashi T, Kikko Y, Kontani K, Katada T. 2008. Domain architecture of the
1320 atypical Arf-family GTPase Arl13b involved in cilia formation. *Biochem Biophys Res*
1321 *Commun.* 373:119–124. doi: 10.1016/j.bbrc.2008.06.001.
- 1322 Houghton FJ et al. 2012. Arl5b is a Golgi-localised small G protein involved in the regulation
1323 of retrograde transport. *Exp Cell Res.* 318:464–477. doi: 10.1016/j.yexcr.2011.12.023.
- 1324 Insall R, Gaudet P, Weeks G. 2005. The small GTPase superfamily. In: Loomis WF, Kuspa
1325 A, editors. *Dictyostelium Genomics*. Norfolk: Horizon Bioscience. p. 173–210.
- 1326 Ishida M, E Oguchi M, Fukuda M. 2016. Multiple Types of Guanine Nucleotide Exchange
1327 Factors (GEFs) for Rab Small GTPases. *Cell Struct Funct.* 41:61–79. doi:
1328 10.1247/csf.16008.
- 1329 Ivanova AA et al. 2017. Biochemical characterization of purified mammalian ARL13B
1330 protein indicates that it is an atypical GTPase and ARL3 guanine nucleotide exchange
1331 factor (GEF). *J Biol Chem.* 292:11091–11108. doi: 10.1074/jbc.M117.784025.
- 1332 Jackson CL. 2014. Arf Proteins and Their Regulators: At the Interface Between Membrane
1333 Lipids and the Protein Trafficking Machinery. In: *Ras Superfamily Small G Proteins:*
1334 *Biology and Mechanisms 2: Transport*. Wittinghofer, A, editor. Springer International
1335 Publishing: Cham pp. 151–180. doi: 10.1007/978-3-319-07761-1_8.
- 1336 Jackson CL, Bouvet S. 2014. Arfs at a glance. *J Cell Sci.* 127:4103–4109. doi:
1337 10.1242/jcs.144899.

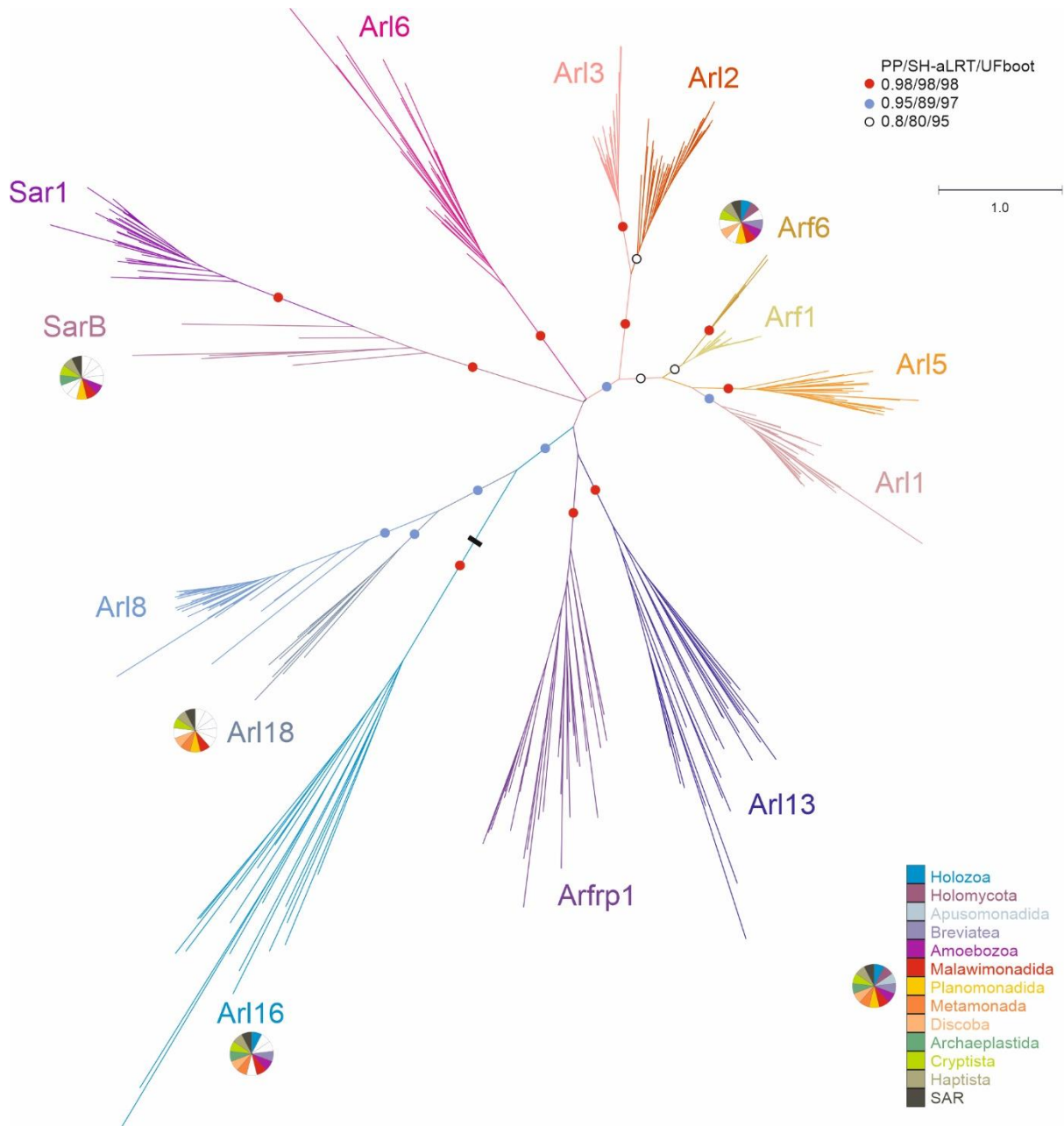
- 1338 Jivan A, Earnest S, Juang Y-C, Cobb MH. 2009. Radial spoke protein 3 is a mammalian
1339 protein kinase A-anchoring protein that binds ERK1/2. *J Biol Chem.* 284:29437–29445.
1340 doi: 10.1074/jbc.M109.048181.
- 1341 Jones S et al. 1999. Genetic interactions in yeast between Ypt GTPases and Arf guanine
1342 nucleotide exchangers. *Genetics.* 152:1543–1556.
- 1343 Kaessmann H, Vinckenbosch N, Long M. 2009. RNA-based gene duplication: mechanistic
1344 and evolutionary insights. *Nat Rev Genet.* 10:19–31. doi: 10.1038/nrg2487.
- 1345 Kahn RA et al. 2008. Consensus nomenclature for the human ArfGAP domain-containing
1346 proteins. *J Cell Biol.* 182:1039–1044. doi: 10.1083/jcb.200806041.
- 1347 Kahn RA et al. 2006. Nomenclature for the human Arf family of GTP-binding proteins: ARF,
1348 ARL, and SAR proteins. *J Cell Biol.* 172:645–650. doi: 10.1083/jcb.200512057.
- 1349 Kahn RA. 2009. Toward a model for Arf GTPases as regulators of traffic at the Golgi. *FEBS*
1350 *Lett.* 583:3872–3879. doi: 10.1016/j.febslet.2009.10.066.
- 1351 Kahn RA, Goddard C, Newkirk M. 1988. Chemical and immunological characterization of
1352 the 21-kDa ADP-ribosylation factor of adenylate cyclase. *J. Biol. Chem.* 263:8282–8287.
- 1353 Katoh K, Standley DM. 2013. MAFFT Multiple Sequence Alignment Software Version 7:
1354 Improvements in Performance and Usability. *Mol Biol Evol.* 30:772–780. doi:
1355 10.1093/molbev/mst010.
- 1356 Keenan RJ, Freymann DM, Stroud RM, Walter P. 2001. The signal recognition particle.
1357 *Annu. Rev. Biochem.* 70:755–775. doi: 10.1146/annurev.biochem.70.1.755.
- 1358 Khatter D, Sindhwani A, Sharma M. 2015. Arf-like GTPase Arl8: Moving from the periphery
1359 to the center of lysosomal biology. *Cell Logist.* 5. doi: 10.1080/21592799.2015.1086501.
- 1360 Klinger CM, et al. 2016. Resolving the homology-function relationship through comparative
1361 genomics of membrane-trafficking machinery and parasite cell biology. *Mol Biochem*
1362 *Parasitol.* 209:88–103. doi: 10.1016/j.molbiopara.2016.07.003.
- 1363 Klinger CM, Spang A, Dacks JB, Ettema TJ. 2016. Tracing the Archaeal Origins of
1364 Eukaryotic Membrane-Trafficking System Building Blocks. *Mol Biol Evol.* 33:1528–
1365 1541. doi: 10.1093/molbev/msw034.
- 1366 Klöpper TH, Kienle N, Fasshauer D, Munro S. 2012. Untangling the evolution of Rab G
1367 proteins: implications of a comprehensive genomic analysis. *BMC Biol.* 10:71. doi:
1368 10.1186/1741-7007-10-71.
- 1369 Kumari B, Kumar R, Kumar M. 2014. PalmPred: an SVM based palmitoylation prediction
1370 method using sequence profile information. *PLoS ONE.* 9:e89246. doi:
1371 10.1371/journal.pone.0089246.
- 1372 Lee FJ et al. 1997. Characterization of an ADP-ribosylation factor-like 1 protein in
1373 *Saccharomyces cerevisiae*. *J Biol Chem.* 272:30998–31005. doi:
1374 10.1074/jbc.272.49.30998.
- 1375 Leipe DD, Wolf YI, Koonin EV, Aravind L. 2002. Classification and evolution of P-loop
1376 GTPases and related ATPases. *J. Mol. Biol.* 317:41–72. doi: 10.1006/jmbi.2001.5378.
- 1377 Lemmon MA. 2007. Pleckstrin homology (PH) domains and phosphoinositides. *Biochem*
1378 *Soc Symp.* 81–93. doi: 10.1042/BSS0740081.
- 1379 Letunic I, Bork P. 2018. 20 years of the SMART protein domain annotation resource.
1380 *Nucleic Acids Res.* 46:D493–D496. doi: 10.1093/nar/gkx922.
- 1381 Li C-C et al. 2007. ARL4D recruits cytohesin-2/ARNO to modulate actin remodeling. *Mol*
1382 *Biol Cell.* 18:4420–4437. doi: 10.1091/mbc.e07-02-0149.
- 1383 Li S et al. 2015. In Silico Identification of Protein S-Palmitoylation Sites and Their
1384 Involvement in Human Inherited Disease. *J Chem Inf Model.* 55:2015–2025. doi:
1385 10.1021/acs.jcim.5b00276.

- 1386 Li Y et al. 2004. Functional genomic analysis of the ADP-ribosylation factor family of
1387 GTPases: phylogeny among diverse eukaryotes and function in *C. elegans*. *FASEB J.*
1388 18:1834–1850. doi: 10.1096/fj.04-2273com.
- 1389 Lin C-Y, Li C-C, Huang P-H, Lee F-JS. 2002. A developmentally regulated ARF-like 5
1390 protein (ARL5), localized to nuclei and nucleoli, interacts with heterochromatin protein 1.
1391 *J Cell Sci.* 115:4433–4445. doi: 10.1242/jcs.00123.
- 1392 Liu Y, Kahn RA, Prestegard JH. 2010. Dynamic structure of membrane-anchored Arf*GTP.
1393 *Nat Struct Mol Biol.* 17:876–881. doi: 10.1038/nsmb.1853.
- 1394 Liu Y, Kahn RA, Prestegard JH. 2009. Structure and Membrane Interaction of Myristoylated
1395 ARF1. *Structure.* 17:79–87. doi: 10.1016/j.str.2008.10.020.
- 1396 Makiuchi T, Nozaki T. 2014. Highly divergent mitochondrion-related organelles in anaerobic
1397 parasitic protozoa. *Biochimie.* 100:3–17. doi: 10.1016/j.biochi.2013.11.018.
- 1398 Manolea F et al. 2010. Arf3 Is Activated Uniquely at the trans-Golgi Network by Brefeldin
1399 A-inhibited Guanine Nucleotide Exchange Factors. *Mol Biol Cell.* 21:1836–1849. doi:
1400 10.1091/mbc.E10-01-0016.
- 1401 Marchler-Bauer A et al. 2015. CDD: NCBI’s conserved domain database. *Nucleic Acids Res.*
1402 43:D222–226. doi: 10.1093/nar/gku1221.
- 1403 Maurer-Stroh S, Eisenhaber B, Eisenhaber F. 2002. N-terminal N-myristoylation of proteins:
1404 prediction of substrate proteins from amino acid sequence. *J. Mol. Biol.* 317:541–557.
1405 doi: 10.1006/jmbi.2002.5426.
- 1406 Melville DB, Studer S, Schekman R. 2020. Small sequence variations between two
1407 mammalian paralogs of the small GTPase SAR1 underlie functional differences in coat
1408 protein complex II assembly. *J Biol Chem.* 295:8401–8412. doi:
1409 10.1074/jbc.RA120.012964.
- 1410 Meza I, Talamás-Rohana P, Vargas MA. 2006. The Cytoskeleton of *Entamoeba histolytica*:
1411 Structure, Function, and Regulation by Signaling Pathways. *Arch Med Res.* 37:234–243.
1412 doi: 10.1016/j.arcmed.2005.09.008.
- 1413 Miertschke M, Koerner C, Spoerner M, Wittinghofer A. 2014. Structural insights into the
1414 small G-protein Arl13B and implications for Joubert syndrome. *Biochem J.* 457:301–311.
1415 doi: 10.1042/BJ20131097.
- 1416 Miller EA, Barlowe C. 2010. Regulation of coat assembly--sorting things out at the ER. *Curr*
1417 *Opin Cell Biol.* 22:447–453. doi: 10.1016/j.ceb.2010.04.003.
- 1418 Miller MA, Pfeiffer W, Schwartz T. 2010. Creating the CIPRES Science Gateway for
1419 inference of large phylogenetic trees. In: 2010 Gateway Computing Environments
1420 Workshop (GCE). pp. 1–8. doi: 10.1109/GCE.2010.5676129.
- 1421 Minh BQ, Nguyen MAT, von Haeseler A. 2013. Ultrafast approximation for phylogenetic
1422 bootstrap. *Mol. Biol. Evol.* 30:1188–1195. doi: 10.1093/molbev/mst024.
- 1423 Mizuno-Yamasaki E, Rivera-Molina F, Novick P. 2012. GTPase networks in membrane
1424 traffic. *Annu Rev Biochem.* 81:637–659. doi: 10.1146/annurev-biochem-052810-093700.
- 1425 More K, Klinger CM, Barlow LD, Dacks JB. 2020. Evolution and Natural History of
1426 Membrane Trafficking in Eukaryotes. *Curr Biol.* 30:R553–R564. doi:
1427 10.1016/j.cub.2020.03.068.
- 1428 Mourão A, Nager AR, Nachury MV, Lorentzen E. 2014. Structural basis for membrane
1429 targeting of the BBSome by ARL6. *Nat Struct Mol Biol.* 21:1035–1041. doi:
1430 10.1038/nsmb.2920.
- 1431 Nabais C, Peneda C, Bettencourt-Dias M. 2020. Evolution of centriole assembly. *Curr Biol.*
1432 30:R494–R502. doi: 10.1016/j.cub.2020.02.036.
- 1433 Neuwald AF. 2007. $G\alpha$ - $G\beta\gamma$ dissociation may be due to retraction of a buried lysine and
1434 disruption of an aromatic cluster by a GTP-sensing Arg Trp pair. *Protein Sci.* 16:2570–
1435 2577. doi: 10.1110/ps.073098107.

- 1436 Newman LE, Schiavon CR, Turn RE, Kahn RA. 2017. The ARL2 GTPase regulates
1437 mitochondrial fusion from the intermembrane space. *Cell Logist.* 7:e1340104. doi:
1438 10.1080/21592799.2017.1340104.
- 1439 Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: a fast and effective
1440 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.*
1441 32:268–274. doi: 10.1093/molbev/msu300.
- 1442 Panic B, Whyte JRC, Munro S. 2003. The ARF-like GTPases Arl1p and Arl3p Act in a
1443 Pathway that Interacts with Vesicle-Tethering Factors at the Golgi Apparatus. *Curr Biol.*
1444 13:405–410. doi: 10.1016/S0960-9822(03)00091-5.
- 1445 Pasqualato S, Renault L, Cherfils J. 2002. Arf, Arl, Arp and Sar proteins: a family of GTP-
1446 binding proteins with a structural device for ‘front–back’ communication. *EMBO Rep.*
1447 3:1035–1041. doi: 10.1093/embo-reports/kvf221.
- 1448 Patel M, Chiang T-C, Tran V, Lee F-JS, Côté J-F. 2011. The Arf family GTPase Arl4A
1449 complexes with ELMO proteins to promote actin cytoskeleton remodeling and reveals a
1450 versatile Ras-binding domain in the ELMO proteins family. *J Biol Chem.* 286:38969–
1451 38979. doi: 10.1074/jbc.M111.274191.
- 1452 Pereira-Leal JB. 2008. The Ypt/Rab family and the evolution of trafficking in fungi. *Traffic.*
1453 9:27–38. doi: 10.1111/j.1600-0854.2007.00667.x.
- 1454 Petrželková R, Eliáš M. 2014. Contrasting patterns in the evolution of the Rab GTPase family
1455 in Archaeplastida. *Acta Soc Bot Polon.* 83:303–315. doi: 10.5586/asbp.2014.052.
- 1456 Pipaliya SV, Schlacht A, Klinger CM, Kahn RA, Dacks J. 2019. Ancient complement and
1457 lineage-specific evolution of the Sec7 ARF GEF proteins in eukaryotes. *Mol Biol Cell.*
1458 30:1846–1863. doi: 10.1091/mbc.E19-01-0073.
- 1459 Resh MD. 1999. Fatty acylation of proteins: new insights into membrane targeting of
1460 myristoylated and palmitoylated proteins. *Biochim Biophys Acta.* 1451:1–16. doi:
1461 10.1016/s0167-4889(99)00075-0.
- 1462 Rojas AM, Fuentes G, Rausell A, Valencia A. 2012. The Ras protein superfamily:
1463 Evolutionary tree and role of conserved amino acids. *J Cell Biol.* 196:189–201. doi:
1464 10.1083/jcb.201103008.
- 1465 Ronquist F et al. 2012. MrBayes 3.2: efficient Bayesian phylogenetic inference and model
1466 choice across a large model space. *Syst Biol.* 61:539–542. doi: 10.1093/sysbio/sys029.
- 1467 Rosa-Ferreira C, Christis C, Torres IL, Munro S. 2015. The small G protein Arl5 contributes
1468 to endosome-to-Golgi traffic by aiding the recruitment of the GARP complex to the
1469 Golgi. *Biol Open.* 4:474–481. doi: 10.1242/bio.201410975.
- 1470 Roy D, Lohia A. 2004. Sequence divergence of *Entamoeba histolytica* tubulin is responsible
1471 for its altered tertiary structure. *Biochem Biophys Res Commun.* 319:1010–1016. doi:
1472 10.1016/j.bbrc.2004.05.079.
- 1473 Roy K et al. 2017. Palmitoylation of the ciliary GTPase ARL13b is necessary for its stability
1474 and its role in cilia formation. *J. Biol. Chem.* 292:17703–17717. doi:
1475 10.1074/jbc.M117.792937.
- 1476 Sarma GN et al. 2010. Structure of D-AKAP2:PKA RI complex: Insights into AKAP
1477 specificity and selectivity. *Structure.* 18:155–166. doi: 10.1016/j.str.2009.12.012.
- 1478 Sato K, Nakano A. 2007. Mechanisms of COPII vesicle formation and protein sorting. *FEBS*
1479 *Lett.* 581:2076–2082. doi: 10.1016/j.febslet.2007.01.091.
- 1480 Schlacht A, Dacks JB. 2015. Unexpected ancient paralogs and an evolutionary model for the
1481 COPII coat complex. *Genome Biol Evol.* 7:1098–1109. doi: 10.1093/gbe/evv045.
- 1482 Schmidt HA, Strimmer K, Vingron M, von Haeseler A. 2002. TREE-PUZZLE: maximum
1483 likelihood phylogenetic analysis using quartets and parallel computing. *Bioinformatics.*
1484 18:502–504. doi: 10.1093/bioinformatics/18.3.502.

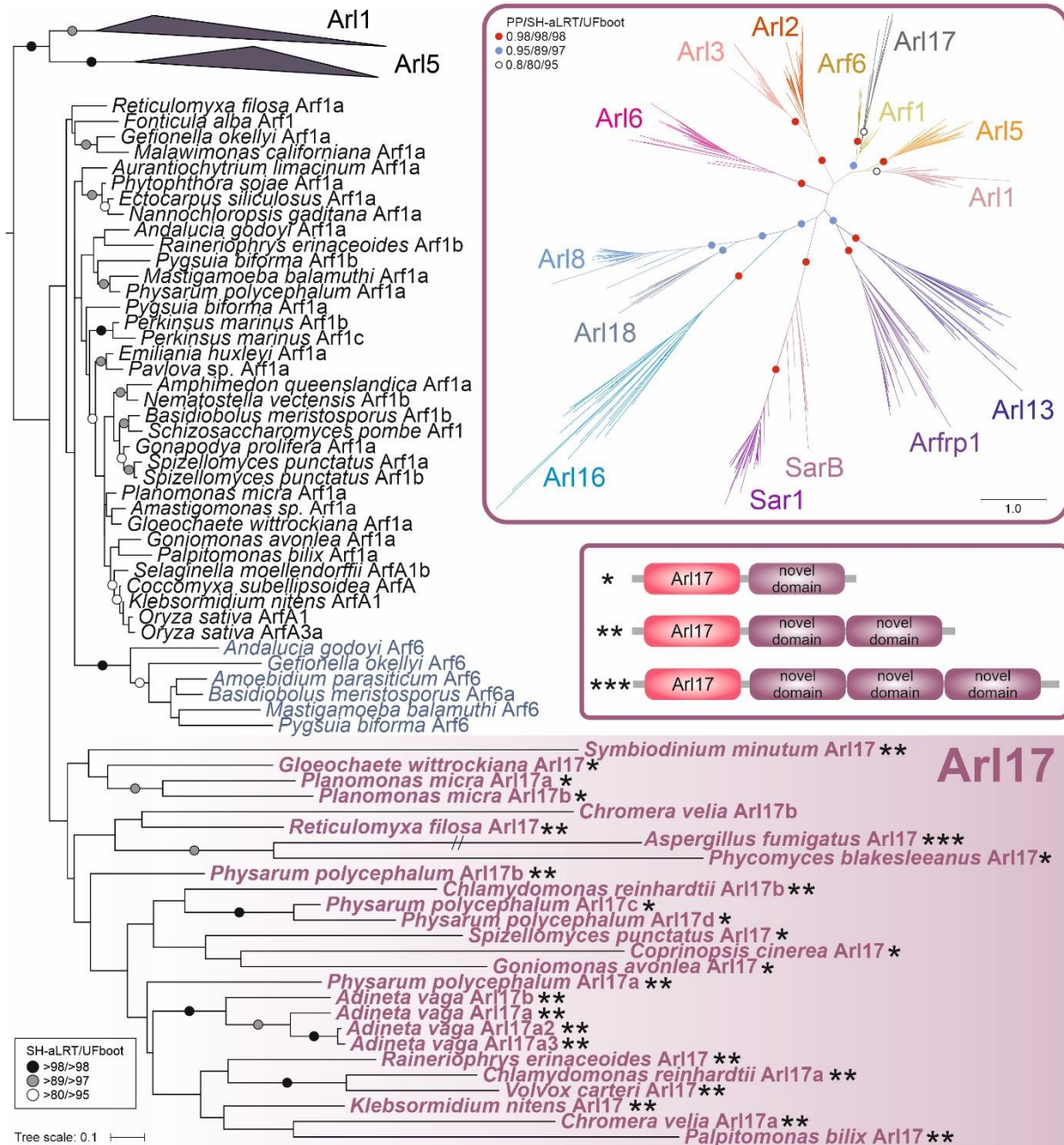
- 1485 Schwartz T, Blobel G. 2003. Structural Basis for the Function of the β Subunit of the
1486 Eukaryotic Signal Recognition Particle Receptor. *Cell*. 112:793–803. doi:
1487 10.1016/S0092-8674(03)00161-2.
- 1488 Schweitzer JK, Sedgwick AE, D’Souza-Schorey C. 2011. ARF6-mediated endocytic
1489 recycling impacts cell movement, cell division and lipid homeostasis. *Semin Cell Dev*
1490 *Biol*. 22:39–47. doi: 10.1016/j.semcdb.2010.09.002.
- 1491 Seidel RD, Amor JC, Kahn RA, Prestegard JH. 2004. Conformational changes in human
1492 Arf1 on nucleotide exchange and deletion of membrane-binding elements. *J Biol Chem*.
1493 279:48307–48318. doi: 10.1074/jbc.M402109200.
- 1494 Setty SRG, Shin ME, Yoshino A, Marks MS, Burd CG. 2003. Golgi Recruitment of GRIP
1495 Domain Proteins by Arf-like GTPase 1 Is Regulated by Arf-like GTPase 3. *Curr Biol*.
1496 13:401–404. doi: 10.1016/S0960-9822(03)00089-7.
- 1497 Setty SRG, Strohlic TI, Tong AHY, Boone C, Burd CG. 2004. Golgi targeting of ARF-like
1498 GTPase Arl3p requires its N α -acetylation and the integral membrane protein Sys1p. *Nat*
1499 *Cell Biol*. 6:414–419. doi: 10.1038/ncb1121.
- 1500 Sharer JD, Shern JF, Van Valkenburgh H, Wallace DC, Kahn RA. 2002. ARL2 and BART
1501 Enter Mitochondria and Bind the Adenine Nucleotide Transporter. *Mol Biol Cell*. 13:71–
1502 83. doi: 10.1091/mbc.01-05-0245.
- 1503 Shi S-P et al. 2013. The prediction of palmitoylation site locations using a multiple feature
1504 extraction method. *J. Mol. Graph. Model*. 40:125–130. doi: 10.1016/j.jmgm.2012.12.006.
- 1505 Smith JC, Northey JGB, Garg J, Pearlman RE, Siu KWM. 2005. Robust Method for
1506 Proteome Analysis by MS/MS Using an Entire Translated Genome: Demonstration on
1507 the Ciliome of *Tetrahymena thermophila*. *J. Proteome Res*. 4:909–919. doi:
1508 10.1021/pr050013h.
- 1509 Söding J, Biegert A, Lupas AN. 2005. The HHpred interactive server for protein homology
1510 detection and structure prediction. *Nucleic Acids Res*. 33:W244–W248. doi:
1511 10.1093/nar/gki408.
- 1512 Stephen LA, Ismail S. 2016. Shuttling and sorting lipid-modified cargo into the cilia.
1513 *Biochem Soc Trans*. 44:1273–1280. doi: 10.1042/BST20160122.
- 1514 Sztul E et al. 2019. ARF GTPases and their GEFs and GAPs: concepts and challenges. *Mol*
1515 *Biol Cell*. 30:1249–1271. doi: 10.1091/mbc.E18-12-0820.
- 1516 Tamkun JW et al. 1991. The arflike gene encodes an essential GTP-binding protein in
1517 *Drosophila*. *Proc Natl Acad Sci U S A*. 88:3120–3124. doi: 10.1073/pnas.88.8.3120.
- 1518 Tria FDK, Landan G, Dagan T. 2017. Phylogenetic rooting using minimal ancestor deviation.
1519 *Nat Ecol Evol*. 1:193. doi: 10.1038/s41559-017-0193.
- 1520 Trifinopoulos J, Nguyen L-T, von Haeseler A, Minh BQ. 2016. W-IQ-TREE: a fast online
1521 phylogenetic tool for maximum likelihood analysis. *Nucleic Acids Res*. 44:W232–W235.
1522 doi: 10.1093/nar/gkw256.
- 1523 Turn RE, East MP, Prekeris R, Kahn RA. 2020. The ARF GAP ELMOD2 acts with different
1524 GTPases to regulate centrosomal microtubule nucleation and cytokinesis. *Mol Biol Cell*.
1525 31:2070–2091. doi: 10.1091/mbc.E20-01-0012.
- 1526 Van Valkenburgh H, Shern JF, Sharer JD, Zhu X, Kahn RA. 2001. ADP-ribosylation factors
1527 (ARFs) and ARF-like 1 (ARL1) have both specific and shared effectors: characterizing
1528 ARL1-binding proteins. *J Biol Chem*. 276:22826–22837. doi: 10.1074/jbc.M102359200.
- 1529 Vernoud V, Horton AC, Yang Z, Nielsen E. 2003. Analysis of the small GTPase gene
1530 superfamily of *Arabidopsis*. *Plant Physiol*. 131:1191–1208. doi: 10.1104/pp.013052.
- 1531 Vetter IR. 2014. The Structure of the G Domain of the Ras Superfamily. In: *Ras Superfamily*
1532 *Small G Proteins: Biology and Mechanisms 1: General Features, Signaling*. Wittinghofer,
1533 A, editor. Springer: Vienna pp. 25–50. doi: 10.1007/978-3-7091-1806-1_2.

- 1534 Vichi A, Moss J, Vaughan M. 2005. ADP-ribosylation factor domain protein 1 (ARD1), a
1535 multifunctional protein with ubiquitin E3 ligase, GAP, and ARF domains. *Meth.*
1536 *Enzymol.* 404:195–206. doi: 10.1016/S0076-6879(05)04019-X.
- 1537 Vlahou G, Eliáš M, von Kleist-Retzow J-C, Wiesner RJ, Rivero F. 2011. The Ras related
1538 GTPase Miro is not required for mitochondrial transport in *Dictyostelium discoideum*.
1539 *Eur J Cell Biol.* 90:342–355. doi: 10.1016/j.ejcb.2010.10.012.
- 1540 Vosseberg J et al. 2021. Timing the origin of eukaryotic cellular complexity with ancient
1541 duplications. *Nat Ecol Evol.* 5:92–100. doi: 10.1038/s41559-020-01320-z.
- 1542 Wang X-B, Wu L-Y, Wang Y-C, Deng N-Y. 2009. Prediction of palmitoylation sites using
1543 the composition of k-spaced amino acid pairs. *Protein Eng. Des. Sel.* 22:707–712. doi:
1544 10.1093/protein/gzp055.
- 1545 Wilson GM et al. 2005. The FIP3-Rab11 protein complex regulates recycling endosome
1546 targeting to the cleavage furrow during late cytokinesis. *Mol Biol Cell.* 16:849–860. doi:
1547 10.1091/mbc.e04-10-0927.
- 1548 Wirschell M et al. 2008. Building a radial spoke: flagellar radial spoke protein 3 (RSP3) is a
1549 dimer. *Cell Motil Cytoskeleton.* 65:238–248. doi: 10.1002/cm.20257.
- 1550 Wuichet K, Søgaard-Andersen L. 2014. Evolution and Diversity of the Ras Superfamily of
1551 Small GTPases in Prokaryotes. *Genome Biol Evol.* 7:57–70. doi: 10.1093/gbe/evu264.
- 1552 Xie Y et al. 2016. GPS-Lipid: a robust tool for the prediction of multiple lipid modification
1553 sites. *Sci Rep.* 6:28249. doi: 10.1038/srep28249.
- 1554 Xu Y, Wang Z, Li C, Chou K-C. 2017. iPreny-PseAAC: Identify C-terminal Cysteine
1555 Prenylation Sites in Proteins by Incorporating Two Tiers of Sequence Couplings into
1556 PseAAC. *Med Chem.* 13:544–551. doi: 10.2174/1573406413666170419150052.
- 1557 Yang Y-K et al. 2011. ARF-like protein 16 (ARL16) inhibits RIG-I by binding with its C-
1558 terminal domain in a GTP-dependent manner. *J Biol Chem.* 286:10568–10580. doi:
1559 10.1074/jbc.M110.206896.
- 1560 Yoon HS et al. 2017. Rhodophyta. In: *Handbook of the Protists*. Archibald, JM, Simpson,
1561 AGB, & Slamovits, CH, editors. Springer International Publishing: Cham pp. 89–133.
1562 doi: 10.1007/978-3-319-28149-0_33.
- 1563 Yu C-J, Lee F-JS. 2017. Multiple activities of Arl1 GTPase in the trans-Golgi network. *J Cell*
1564 *Sci.* 130:1691–1699. doi: 10.1242/jcs.201319.
- 1565 Zahn C et al. 2006. Knockout of Arfrp1 leads to disruption of ARF-like1 (ARL1) targeting to
1566 the trans-Golgi in mouse embryos and HeLa cells. *Mol Membr Biol.* 23:475–485. doi:
1567 10.1080/09687860600840100.
- 1568 Záhonová K et al. 2018. Extensive molecular tinkering in the evolution of the membrane
1569 attachment mode of the Rheb GTPase. *Sci Rep.* 8:5239. doi: 10.1038/s41598-018-23575-
1570 0.
- 1571 Zhang Y et al. 2018. The unusual flagellar-targeting mechanism and functions of the
1572 trypanosome ortholog of the ciliary GTPase Arl13b. *J Cell Sci.* 131. doi:
1573 10.1242/jcs.219071.
- 1574 Zhou B, Cox AD. 2014. Posttranslational Modifications of Small G Proteins. In: *Ras*
1575 *Superfamily Small G Proteins: Biology and Mechanisms 1: General Features, Signaling*.
1576 Wittinghofer, A, editor. Springer: Vienna pp. 99–131. doi: 10.1007/978-3-7091-1806-
1577 1_5.



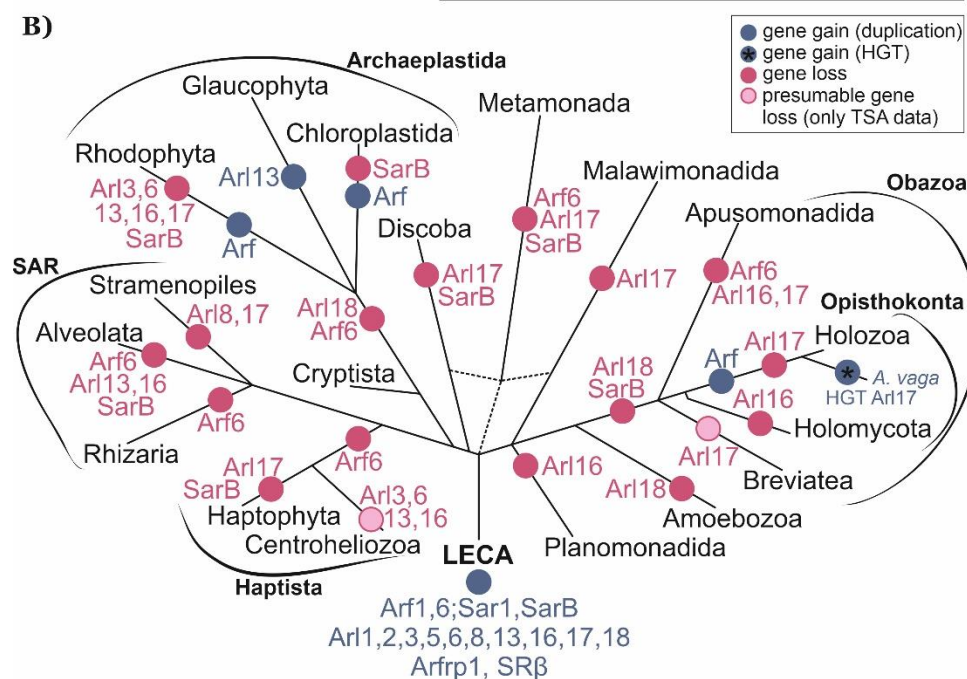
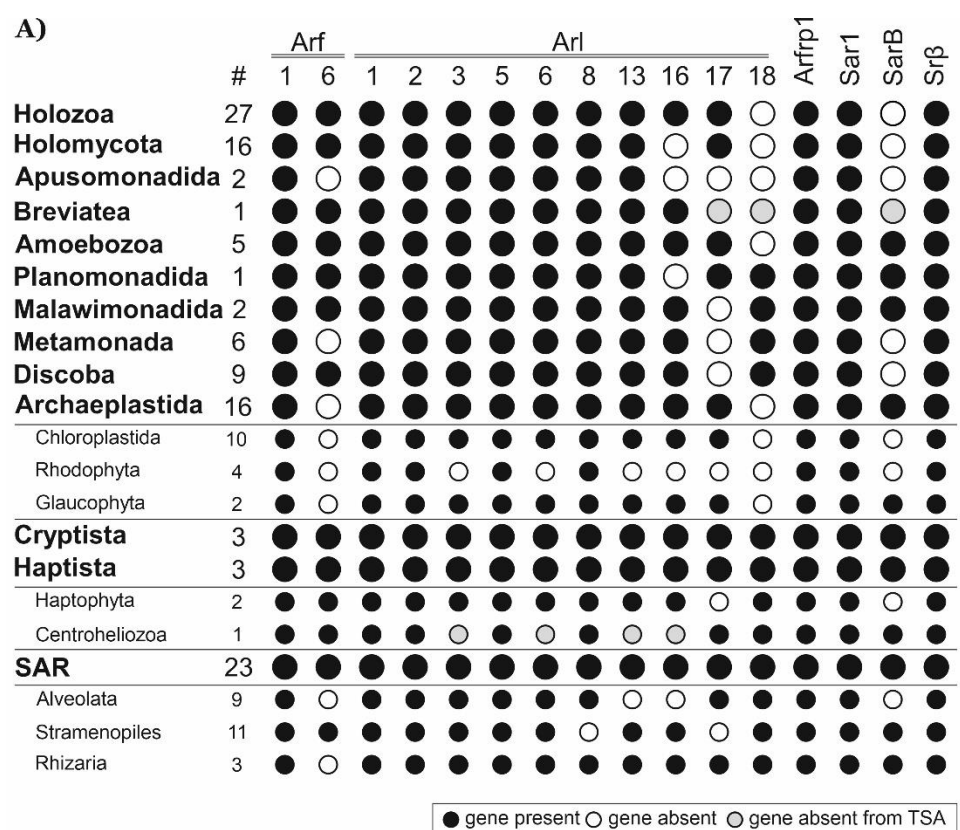
1578
1579

1580 **Fig. 1. Maximum likelihood phylogenetic tree of the ARF family based on a reduced ScrollSaw**
1581 **dataset.** The tree was inferred using IQ-TREE with LG+I+G4 model (the model selected by the
1582 program itself) based on a multiple alignment of 348 protein sequences. Branch support was evaluated
1583 with MrBayes (posterior probability, PP) and with IQ-TREE using the SH-aLRT test and the ultrafast
1584 (UF) bootstrap algorithm (both 10,000 replicates), as described under Materials and Methods. Dots at
1585 branches represent bootstrap values as indicated in the graphical legend (top right), the black bar
1586 indicates the position of the root of the tree as determined with the MAD method. The bar on the top
1587 corresponds to the estimated number of substitutions per site. The pie charts indicate the occurrence of
1588 Arf6, Arl16, Arl18 and SarB in main eukaryotic lineages (indicated by different colours explained in the
1589 graphical legend in the lower right). The remaining paralogs have ubiquitous distribution (i.e., are
1590 present in all main lineages analysed). A full version of the tree is provided in supplementary fig. 2,
1591 Supplementary Material online.



1592
1593

1594 **Fig. 2. Phylogenetic analysis and domain architecture of Arl17.** The tree shown is a result of a ML
 1595 analysis of all Arl17 sequences and a subset of the reduced “scrollsawed” dataset restricted to Arf1,
 1596 Arf6, Arl1, and Arl5 sequences (the latter two collapsed as triangles), altogether 127 protein sequences.
 1597 The alignment was trimmed manually. The tree was inferred using IQ-TREE with LG+I+G4 model (the
 1598 model selected by the program itself) with the ultrafast bootstrap algorithm and the SH-aLRT test (both
 1599 10,000 replicates). Dots at branches represent bootstrap values as indicated in the graphical legend
 1600 (top right). The upper inset shows the ML tree inferred from a full reduced “scrollsawed” dataset
 1601 combined with a subset of Arl17 sequences (picking one representative per each major eukaryote
 1602 group), altogether 356 protein sequences. The tree was inferred using the same approach as the tree
 1603 shown in fig. 1. The inset beneath provides a schematic representation of three different variants of the
 1604 Arl17 domain architecture (correspondence to specific proteins in the tree is indicated by the asterisks).
 1605 The exact architecture of the *Ch. velia* Arl17b protein could not be determine due to incompleteness of
 1606 the genome assembly.



1607

1608

1609

1610

1611

1612

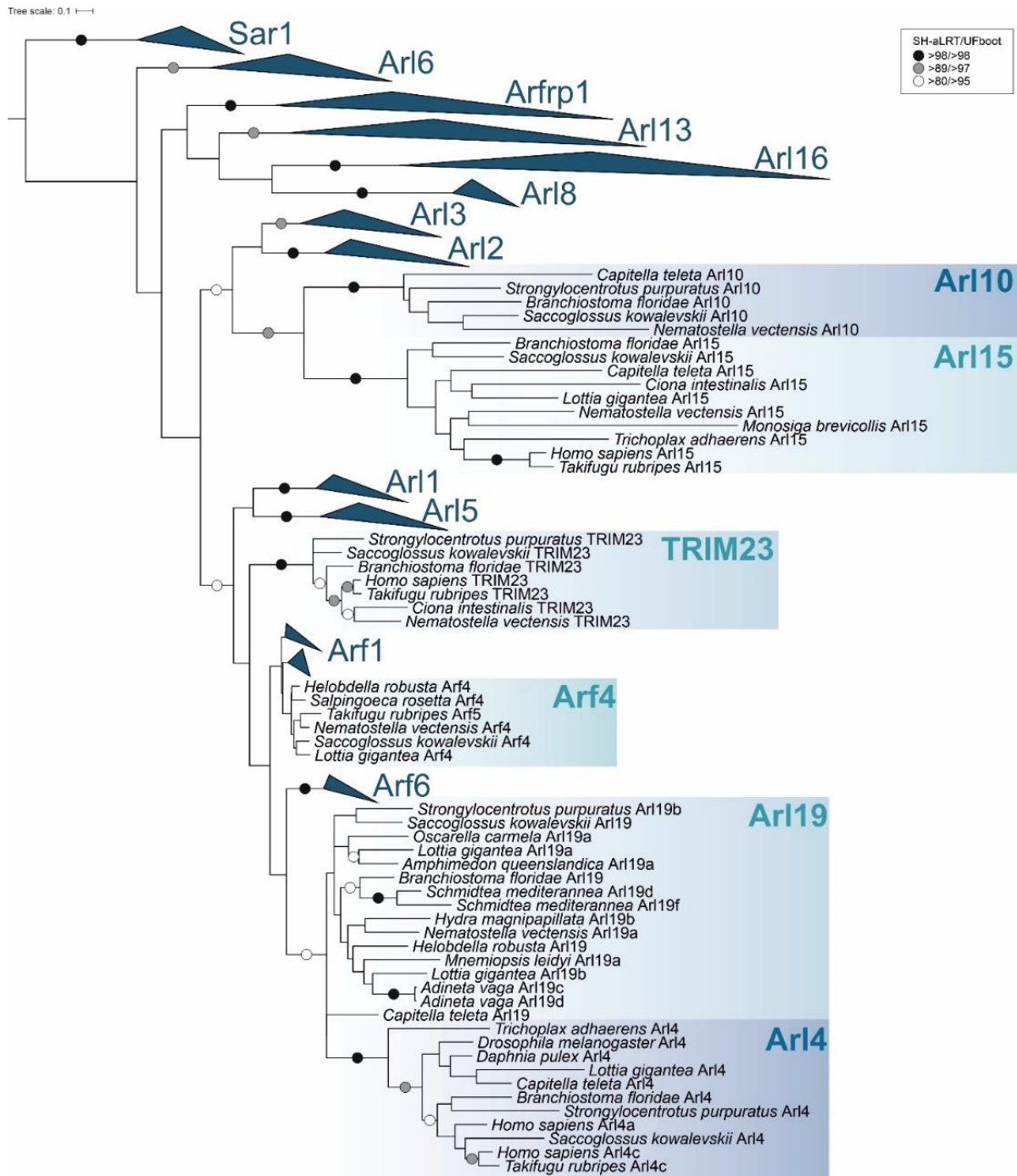
1613

1614

1615

1616

Fig. 3. Retention of ancient paralogs of the ARF family in main lineages of eukaryotes. (A) Black circle: the paralog is present in at least one member of the lineage. White circle: the gene is absent from the lineage (evidenced by genome sequence data). Grey circle: the gene was not found in the transcriptome data available (lineages with transcriptome assemblies only). The hashtag (#) indicates the number of species included in the analysis. **(B)** Gene gains (blue circles) and losses (pink circles) mapped onto the eukaryote phylogeny. Only duplications specific to whole lineages listed in the picture are considered. The acquisition of Arl17 via HGT in rotifers (here represented by *A. vaga*) is indicated with a blue circle with an asterisk within.



1617

1618

1619 **Fig. 4. Maximum likelihood phylogenetic tree of the ARF family based on a ScrollSaw dataset in**

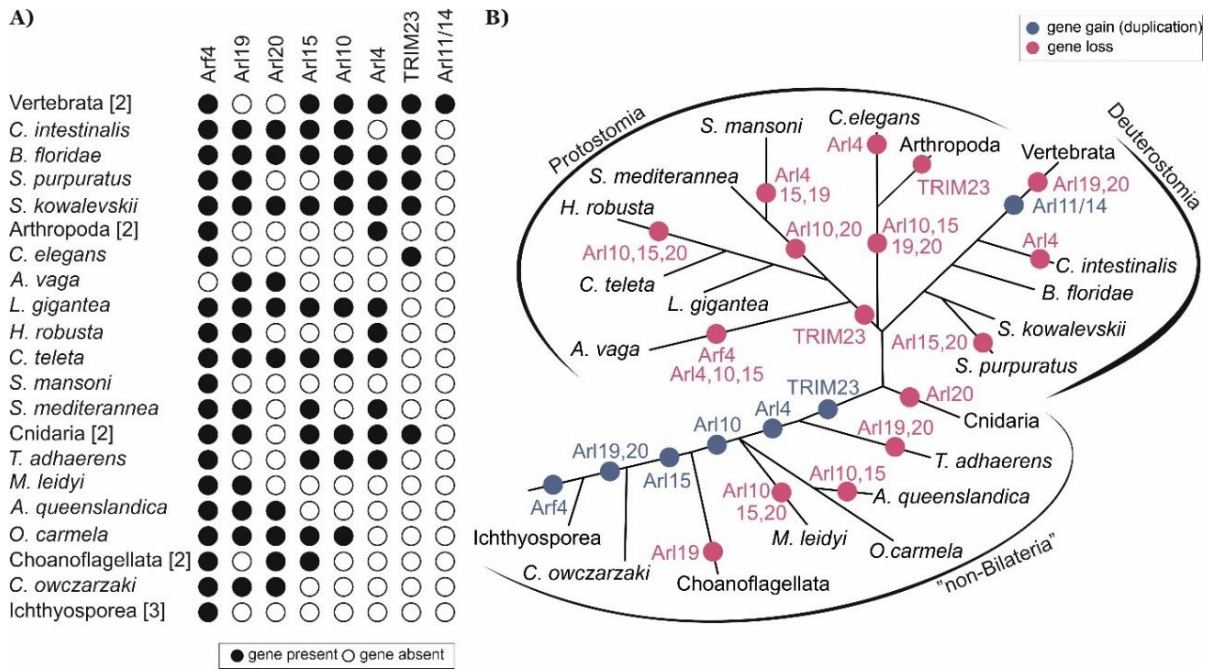
1620 **Holozoa.** The tree was inferred using IQ-TREE with LG++G4 model (the model selected by the

1621 program itself) from a multiple alignment of 323 protein sequences with the ultrafast bootstrap algorithm

1622 and the SH-aLRT test (both 10000 replicates), as described under Materials and Methods. Dots at

1623 branches represent bootstrap values as indicated in the legend shown in the bottom left. Eukaryotic

1624 ancestral paralogs are collapsed as triangles.



1625

1626

1627

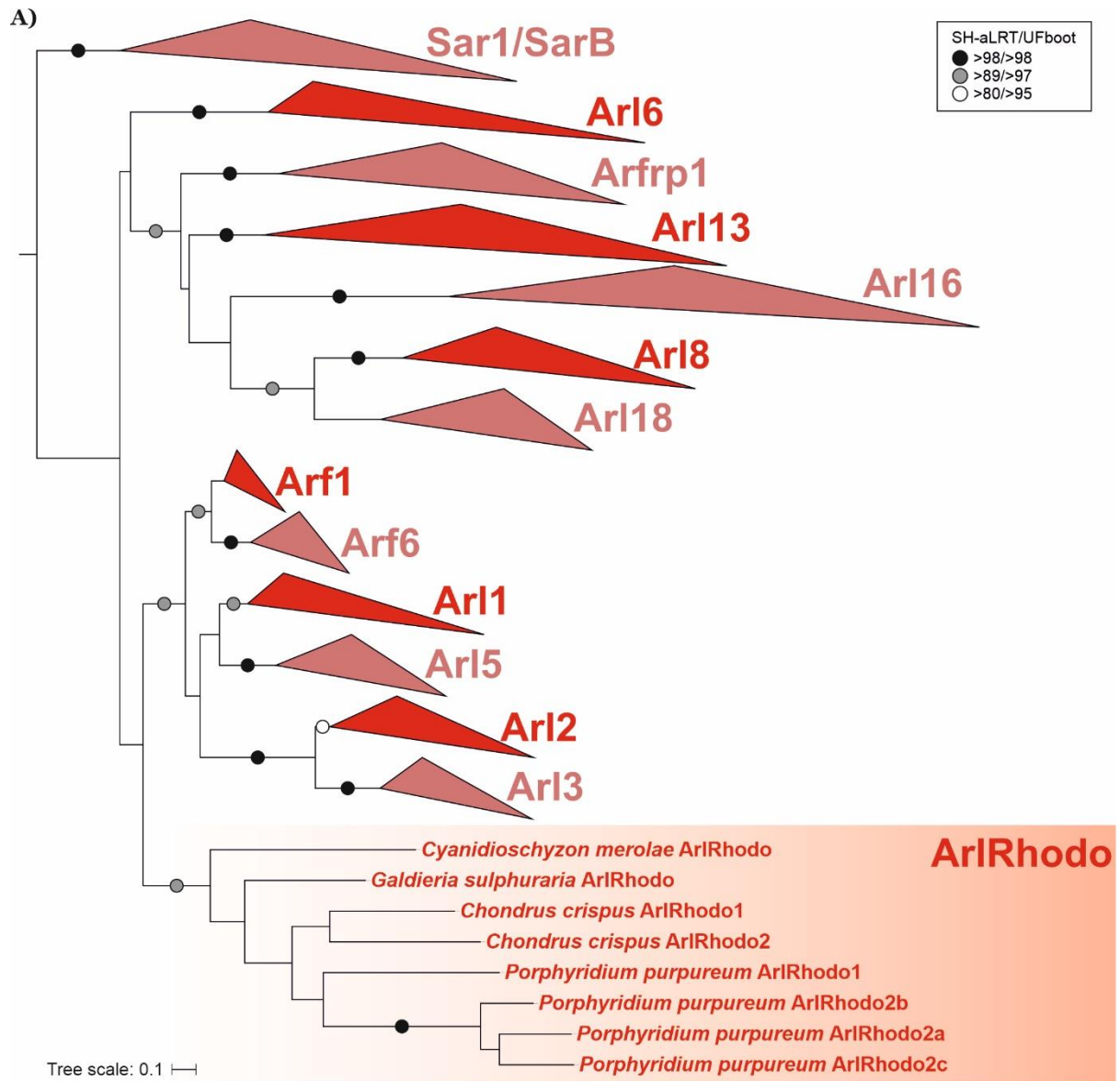
1628

1629

1630

1631

Fig. 5. Retention of lineage-specific paralogs of the ARF family in main lineages of Holozoa. (A) Black circle: the paralog is present in at least one member of the lineage; white circle: the gene is absent from the lineage (evidenced by genome sequence data). Species with identical distribution are collapsed into higher taxa with the number of species indicated in the square brackets. **(B)** Gene gains (blue circles) and losses (pink circles) mapped onto the holozoan phylogeny.



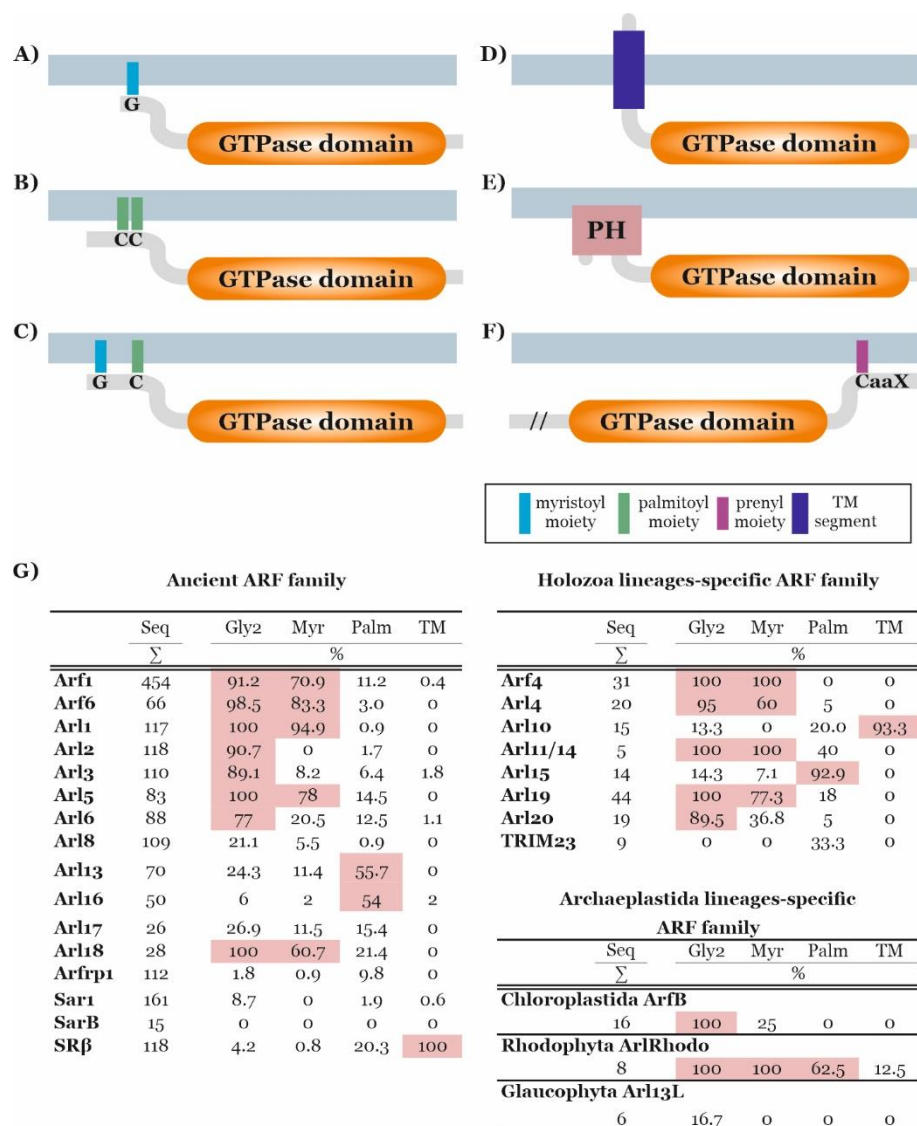
B)

<<<< GTPase domain >>>>

<i>C. merolae</i> ArlRhodo	MGACFGKPSE (4) ARDASSTRRRGSRGD (16) RPGQRIELLLLGLDGAGATTILYRLKLQKF
<i>G. sulphuraria</i> ArlRhodo	MG-CFS-----SKEFQSNQKKK-----YGVVLIIGLDGSGATTILYFLKLGKQ
<i>P. purpureum</i> ArlRhodo1	MG-CFS-----SKEAGARGRGGSKS-----VLFGLDGSSTTVMYQLVLRKQ
<i>P. purpureum</i> ArlRhodo2a	MGQGMSRVYV---EKYKQKNARTGGKG-----ILLIGLDGAGKTTVCYQFLLGKH
<i>P. purpureum</i> ArlRhodo2b	MGGALSTKYV---DNYRARNQKSGRRG-----ILVIGLDGSGKTTVTYQITLQKH
<i>P. purpureum</i> ArlRhodo2c	MGGAVSKAYV---DKYAARNAKTGKRG-----LLVGLDGSSTTVMYQLVLRKQ
<i>C. crispus</i> ArlRhodo1	MGSALS-----ACTGPRGGKDG-----LPC-----VLLVGLDGSSTTILYQVKGKQ
<i>C. crispus</i> ArlRhodo2	MGACMS-----ASHFRHAGSRG-----VLVGLDGSSTTILYQLVLRKQ
Consensus/80%	MG.sbS.....p..psGp+t.....lLlIGLDGtG.TTlhYpb.LGKp

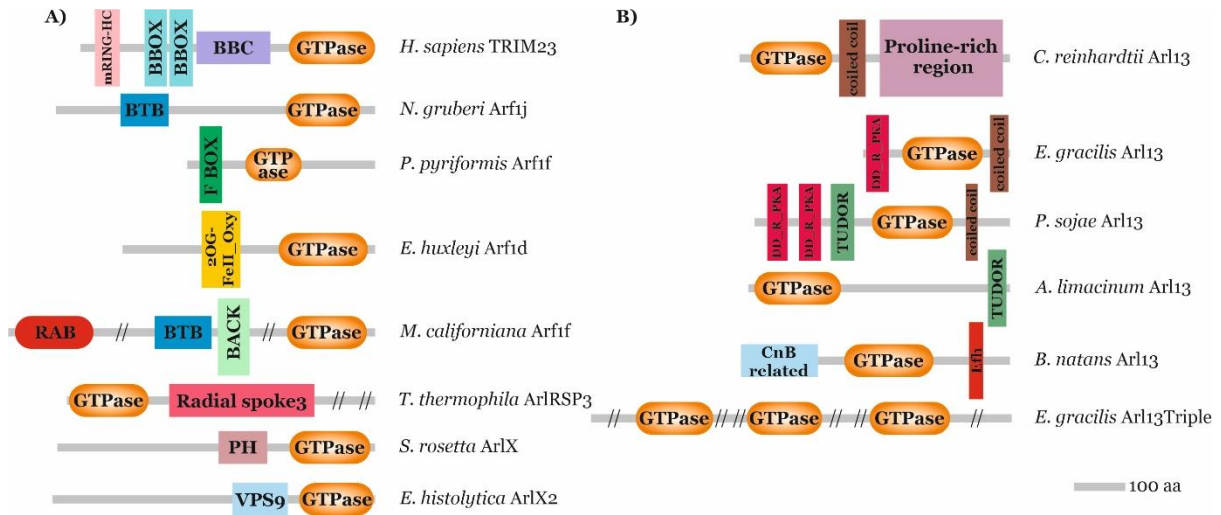
1632
1633

1634 **Fig. 6. ArlRhodo, a novel ARF family member specific for red algae.** (A) Phylogenetic analysis of
 1635 the ArlRhodo group. The tree shown is a result of a ML analysis of all ArlRhodo sequences and the
 1636 reduced “scrollsawed” dataset (altogether 356 sequences). The tree was inferred using IQ-TREE with
 1637 LG+I+G4 model (the model selected by the program itself) with the ultrafast bootstrap algorithm and the
 1638 SH-aLRT test (both 10,000 replicates). Dots at branches represent bootstrap values as indicated in the
 1639 graphical legend (top right). (B) N-terminal region of ArlRhodo proteins with the characteristic
 1640 configuration of glycine and cysteine residues (highlighted in red) predicted to be N-myristoylated and
 1641 S-palmitoylated, respectively.



1642
1643

1644 **Fig. 7. Membrane attachment mechanisms of ARF family proteins.** Examples of different broadly
 1645 conserved mechanisms of membrane attachments of ARF family members are depicted. (A) N-
 1646 terminally myristoylated glycine residues, common for Arfs and several Arf-like proteins. (B) One or two
 1647 S-palmitoylated cysteine residues near the N-terminus, typical for Arl16 and also common in Arl13. (C)
 1648 N-terminally myristoylated glycine residue coupled with S-palmitoylated cysteine residue near the N-
 1649 terminus, typical for ArlRhodo. (D) N-terminal transmembrane region, typical for SRβ and Arl10. (E) N-
 1650 terminally accreted PH domain, present in divergent Arf-like proteins in kinetoplastids and
 1651 choanoflagellates. (F) Prenylation motif (CaaX) at the C-terminus of certain eustigmatophyte-specific
 1652 ARF family members (characterized also by a long N-terminal extension, in the figure marked with "//").
 1653 supplementary table 1, Supplementary Material online lists all identified ARF family proteins predicted
 1654 to be N-myristoylated or S-palmitoylated, or to contain a transmembrane region or PH domain. (G)
 1655 Summary of the results of prediction of N-myristoylation, S-palmitoylation and presence of the
 1656 transmembrane (TM) region in particular subgroups of the ARF family. For each subgroup (group of
 1657 orthologs), the number of sequences (Seq) and the percentages of sequences with glycine residues at
 1658 the second position (Gly2), sequences predicted as N-myristoylated (Myr), sequences predicted as S-
 1659 palmitoylated on at least one cysteine residue (Palm), and sequences with predicted transmembrane
 1660 region(s) (TM) are given. Values above 50% are highlighted in pink. For complete data see
 1661 supplementary tables 1 and 5, Supplementary Material online. These predictions were done as
 1662 described under Materials and Methods.



1663
1664

1665 **Fig. 8. Multi-domain architectures of ARF family proteins.** (A) Examples of lineage-specific ARF
1666 family proteins with extra domains accreted to the GTPase domain. Sequence IDs of the proteins listed
1667 are provided in supplementary table 1, Supplementary Material online. (B) Variation in the domain
1668 architecture of Arl13 proteins across the eukaryote diversity. The Arl13 from *Chlamydomonas*
1669 *reinhardtii* represents the most common and presumably ancestral state.