

1 **Phage origin of mitochondrion-localized family A DNA**
2 **polymerases in kinetoplastids and diplomids**

3

4 Ryo Harada¹ & Yuji Inagaki^{1,2,*}

5

6 ¹Graduate School of Life and Environmental Sciences, University of Tsukuba, Japan.

7 ²Center for Computational Sciences, University of Tsukuba, Japan.

8

9 *Author for Correspondence: Yuji Inagaki, Center for Computational Sciences,
10 University of Tsukuba, Tsukuba, Japan. Phone: +81 29 853 6483, Fax: +81 29 853
11 6404, Email: yuji@ccs.tsukuba.ac.jp

12 Abstract

13 Mitochondria retain their own genomes as other bacterial endosymbiont-derived
14 organelles. Nevertheless, no protein for DNA replication and repair is encoded in any
15 mitochondrial genomes (mtDNAs) assessed to date, suggesting the nucleus
16 primarily governs the maintenance of mtDNA. As the proteins of diverse evolutionary
17 origins occupy a large proportion of the current mitochondrial proteomes, we
18 anticipate finding the same evolutionary trend in the nucleus-encoded machinery for
19 mtDNA maintenance. Indeed, none of the DNA polymerases (DNAPs) in the
20 mitochondrial endosymbiont, a putative α -proteobacterium, seemingly had been
21 inherited by their descendants (mitochondria), as none of the known types of
22 mitochondrion-localized DNAP showed a specific affinity to the α -proteobacterial
23 DNAPs. Nevertheless, we currently have no concrete idea of how and when the
24 known types of mitochondrion-localized DNAPs emerged. We here explored the
25 origins of mitochondrion-localized DNAPs after the improvement of the samplings of
26 DNAPs from bacteria and phages/viruses. Past studies revealed that a set of
27 mitochondrion-localized DNAPs in kinetoplastids and diplomonids, namely PolIB,
28 PolIC, PolID, PolI-Perk1/2, and PolI-dipl (henceforth designated collectively as
29 “PolIBCD+”) have emerged from a single DNAP. In this study, we recovered an
30 intimate connection between PolIBCD+ and the DNAPs found in a particular group of
31 phages. Thus, the common ancestor of kinetoplastids and diplomonids most likely
32 converted a laterally acquired phage DNAP into a mitochondrion-localized DNAP
33 that was ancestral to PolIBCD+. The phage origin of PolIBCD+ hints at a potentially
34 large contribution of proteins acquired via non-vertical processes to the machinery
35 for mtDNA maintenance in kinetoplastids and diplomonids.

36

37 *Keywords:* DNA replication, DNA repair, autographivirus, Euglenozoa, lateral gene
38 transfer, mitochondria

39 Introduction

40 Mitochondria in the extant eukaryotes are the descendants of an
41 endosymbiotic α -proteobacterium in the last eukaryotic common ancestor (Roger et
42 al. 2017). The mitochondrial (mt) proteins, which are localized in mitochondria, are

43 almost entirely nucleus-encoded and evolutionarily multifarious (Gabaldón and
44 Huynen 2007; Wang and Wu 2014; Gray 2015). Only 10-20% of mt proteins were
45 predicted to be of the α -proteobacterial origin, suggesting that the original proteome
46 of the mitochondrial endosymbiont has been remodeled largely (Gray 2015). There
47 are three possible evolutionary paths that coopt non- α -proteobacterial proteins into
48 the molecular machinery in mitochondria. Non- α -proteobacterial mt proteins could
49 emerge (i) de novo, (ii) by recycling of the pre-existing eukaryotic proteins, or (iii) via
50 lateral gene transfer. Mitochondria, in principle, retain their own genomes that have
51 been descended from the mitochondrial endosymbiont, albeit the entire set of
52 proteins required for mtDNA maintenance (replication and repair) is nucleus-
53 encoded. Thus, as a part of the mitochondrial proteome, the machinery for mtDNA
54 maintenance may be dominated by non- α -proteobacterial proteins. Indeed, none of
55 the known DNA polymerases (DNAPs) localized in mitochondria is most unlikely the
56 direct descendant of the DNAPs in the α -proteobacterial endosymbiont that gave rise
57 to the ancestral mitochondrion (see below).

58 Phylogenetically diverse eukaryotes possess family A (famA) DNAPs that are
59 evolutionarily related to DNA polymerase I (Poll) in bacteria (Jung et al. 1987;
60 Moriyama et al. 2011). Some of famA DNAPs in eukaryotes are known to be
61 localized in mitochondria (Krasich and Copeland 2017). So far, four distinct types of
62 mitochondrion-localized famA DNAP have been identified. First, “plant and protist
63 organellar DNA polymerase (POP)” appeared to be broadly distributed among
64 eukaryotes (Moriyama et al. 2011; Hirakawa and Watanabe 2019). Second, animals
65 and fungi are known to use DNA polymerase gamma (Poly) for mtDNA maintenance
66 (Graziewicz et al. 2006). The third type of mitochondrion-localized famA DNAP is
67 “PollA” shared among members of the classes Kinetoplastea, Diplonemea, and
68 Euglenida, which comprise the phylum Euglenozoa (Klingbeil et al. 2002; Harada et
69 al. 2020). Members of Kinetoplastea and Diplonemea possess the fourth type of
70 mitochondrion-localized famA DNAP. “PollB,” “PollC,” and “PollD” were reported
71 originally from a model kinetoplastid *Trypanosoma brucei*, and later identified in
72 broad members of Kinetoplastea (Klingbeil et al. 2002; Harada et al. 2020). The
73 three DNAPs were shown to be closely related to one another in phylogenetic
74 analyses. A recent study further identified multiple DNAPs, which are closely related
75 to but distinct from PollB, C, or D, in an early-branching kinetoplastid *Perkinsella* sp.
76 and diverse diplomonads (Poll-Perk1/2 and Poll-dipl; Harada et al. 2020). PollB, C, D,

77 and their related DNAPs were derived from a single molecule, and thus can be
78 regarded collectively as the fourth type of mitochondrion-localized famA DNAP
79 (henceforth termed as “PolIBCD+” in this study). Pioneering studies considered none
80 of the known mitochondrion-localized famA DNAPs as the direct descendant of PolI
81 in the mitochondrial endosymbiont, but failed to clarify how and when POP, Poly,
82 PolIA, and PolIBCD+ were established in eukaryotic evolution (Moriyama et al. 2011;
83 Hirakawa and Watanabe 2019; Harada et al. 2020).

84 In this study, we explored the origins of mitochondrion-localized famA DNAPs
85 by analyzing an improved dataset wherein sequence sampling from bacteria and
86 phages was improved drastically. We recovered the intimate affinity between
87 PolIBCD+ and the famA DNAPs of a particular group of phages in phylogenetic
88 analyses. Furthermore, these DNAPs appeared to share a unique insertion of
89 consecutive 8 amino acid (aa) residues. Altogether, we conclude that the extent
90 DNAPs belonging to PolIBCD+ were derived from a single phage famA DNAP
91 acquired by the common ancestor of Kinetoplastea and Diplonemea. We also
92 propose that PolIA in Euglenozoa emerged from a type of cytosolic famA DNAP
93 (Pol θ). The origins of PolIA and PolIBCD+ maybe a tip of the remodeling of the
94 machinery of mtDNA maintenance undergone in Kinetoplastea and Diplonemea.

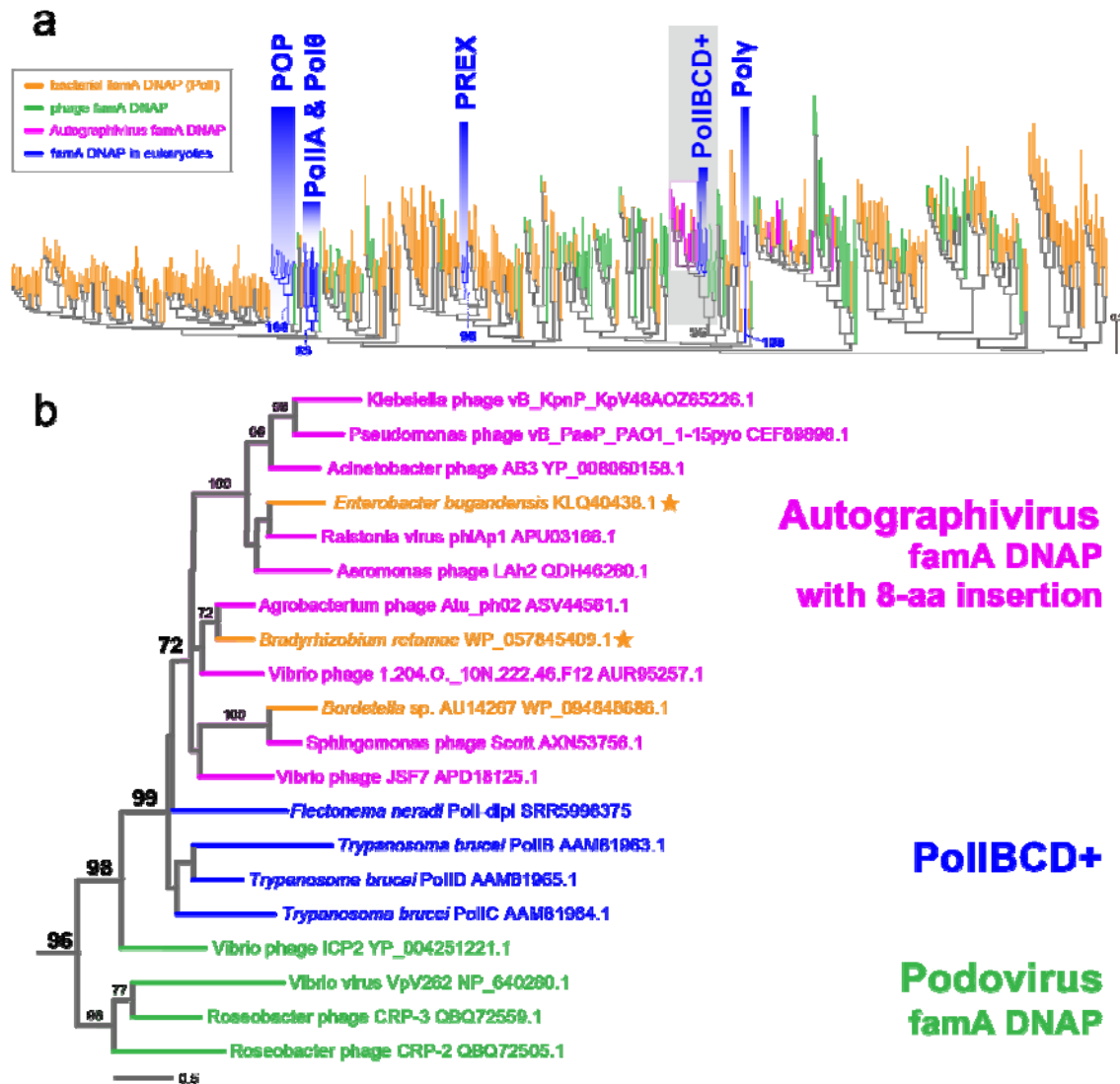
95 Results

96 Prior to this study, the origin of none of the four types of mitochondrion-
97 localized famA DNAPs (i.e. Poly, POP, PolIA, and PolIBCD+) has been elucidated in
98 detail. This study successfully clarified the origin of PolIBCD+ by analyzing
99 phylogenetic alignments that are much richer in bacterial and phage famA DNAPs
100 than those analyzed in the past studies. The sampling of the bacterial homologs was
101 insufficient to reflect the diversity of bacteria in the previously published phylogenies
102 of famA DNAPs (Moriyama et al. 2011; Hirakawa and Watanabe 2019; Harada et al.
103 2020). Furthermore, only a few famA DNAPs of phages have been included in the
104 phylogenetic analyses. In this study we prepared the “global famA DNAP” alignment
105 by incorporating diverse bacterial and phage sequences (446 in total) deposited in
106 public databases and 27 sequences that represent the four mitochondrion-localized
107 types of DNAPs (Poly, POP, PolIA, and PolIBCD+), a single cytosolic DNAP (Pol θ),

108 and a single plastid-localized DNAP found exclusively in apicomplexans and
109 chrompodellids (PREX).

110 The global famA DNAP phylogeny reconstructed four clades, all comprising
111 the eukaryotic homologs exclusively: (i) POP, (ii) PolIA plus Pol θ , (iii) PREX, and (iv)
112 Poly (Shaded in blue in Fig. 1A; see the supplementary materials for the tree with
113 sequence names). The maximum likelihood bootstrap values (MLBPs) for the four
114 clades varied between from 69 to 100%. The POP, PolIA plus Pol θ , or Poly
115 sequences showed no clear affinity to any bacterial or phage famA DNAPs, leaving
116 their origins uncertain. The PREX sequences grouped with bifunctional 3'-5'
117 exonuclease/DNA polymerases in phylogenetically limited bacteria as previously
118 reported (Janouškovec et al. 2015; Hirakawa and Watanabe 2019; Harada et al.
119 2020). Curiously, the PolI/BCD+ sequences were paraphyletic but nested within a
120 robustly supported clade mainly comprising famA DNAP homologs of phages
121 belonging to families Autographiviridae and Podoviridae (Fig. 1B; this figure
122 corresponds to the portion shaded in gray in Fig. 1A). The famA DNAP homologs of
123 autographiviruses and three bacteria formed a subclade with an MLBP of 72%. The
124 coding regions of two out of the three bacterial famA DNAP homologs in this
125 subclade (marked by stars in Fig. 1B) are flanked by phage-like open reading frames
126 (ORFs) in the corresponding genome assemblies deposited under the GenBank
127 accession Nos LEDQ01000001.1 and NZ_LLYA01000167.1. Phage-like ORFs
128 including that of famA DNAP encompass >40 Kbp consecutively in the two bacterial
129 genomes. Thus, the two “bacterial famA DNAPs” are most likely of lysogenic
130 autographiviruses in bacterial genomes. On the other hand, no phage-like ORF was
131 found around that of famA DNAP in the genome of *Bordetella* genomsp. 9 strain
132 AU14267 (NZ_CP021109.1), suggesting that this bacterium acquired a famA DNAP
133 gene from an autographivirus horizontally. The four PolI/BCD+ sequences were
134 positioned at the base of the Autographiviridae clade described above and the
135 grouping of PolI/BCD+ sequences and autographivirus famA DNAPs as a whole
136 received an MLBP of 99% (Fig. 1B). The global famA DNAP phylogeny strongly
137 suggests an intimate evolutionary affinity between PolI/BCD+ and autographivirus
138 famA DNAPs.

139



140

141

142 Fig. 1. Maximum likelihood (ML) phylogenetic tree inferred from an alignment of the famA

143 DNAP sequences of bacteria, phages/viruses, and eukaryotes. (A) Overview of the entire ML

144 tree. All of the sequence names are omitted. The bacterial and eukaryotic sequences are shown

145 in orange and blue, respectively. The sequences of autographiviruses are shown in magenta. A

146 subset of autographiviruses possess famA DNAPs in the pink-shaded clade bears the

147 characteristic insertion of 8 amino acid residues (AGV^{+ins} famA DNAPs; see the main text for the

148 details). Other phage/viral sequences are shown in green. Only ML bootstrap values of interest

149 are shown. The subtree containing PolIBCD+ and AGV^{+ins} famA DNAP sequences (shaded in

150 gray) is enlarged and presented as (B). ML bootstrap values greater than 70% are shown.

151 AGV^{+ins} famA DNAP sequences marked by stars are of the putative lysogenic phages in bacterial

152 genomes.

153

154 Members of Autographiviridae commonly display head-to-tail capsid

155 structures and possess double-stranded linear DNA genomes of approximately 41

156 Kbp in length. This viral family comprises 9 subfamilies and 132 genera (Lavigne et

al. 2008; Adriaenssens et al. 2020). We searched for autographivirus famA DNAPs

157 in the GenBank nr database and detected 175 homologs of 99 members belonging
158 to 57 genera and 76 unclassified members. Each of the 175 members of
159 Autographiviridae seemingly possesses a single famA DNAP. Intriguingly, the
160 autographivirus famA DNAPs were split into two types based on the
161 presence/absence of “8-aa insertion” in the polymerase domain (Fig. S1 and Table
162 S1). In this study, we designate autographivirus famA DNAPs with 8-aa insertion as
163 “AGV^{+ins} famA DNAPs”. Each AGV^{+ins} famA DNAPs was predicted to possess only
164 polymerase domain by InterProScan5 with the Pfam database (Jones et al. 2014; El-
165 Gebali et al. 2019) (Table. S2). AGV^{+ins} famA DNAPs were found in 40 members
166 belonging to 23 genera, and 51 unclassified members (Fig. S1). Although only a
167 subset of the 175 autographivirus famA DNAPs was included, the global famA DNAP
168 phylogeny (Fig. 1A) demonstrated the distant relationship between AGV^{+ins} famA
169 DNAPs and other autographivirus famA DNAPs lacking 8-aa insertions.

170 To reexamine the phylogenetic affinity between PolIBCD+ and AGV^{+ins} famA
171 DNAPs, we selected non-redundant sequences from the 91 AGV^{+ins} famA DNAPs
172 and aligned with 24 PolIBCD+ sequences and four famA DNAPs of phages
173 belonging to a family Podoviridae as the outgroup. The second famA DNAP
174 alignment was subjected to both ML and Bayesian methods. In the second
175 phylogenetic analyses, AGV^{+ins} famA DNAPs and PolIBCD+ sequences formed a
176 clade supported by an MLBP of 100% and a Bayesian posterior probability (BPP) of
177 1.0 (Fig. 2). PolIBCD+ sequences appeared to possess 8 amino acids that are most
178 likely homologous to 8-aa insertion in AGV^{+ins} famA DNAPs (Fig. 2), strengthening
179 the phylogenetic affinity between PolIBCD+ and AGV^{+ins} famA DNAPs. Besides
180 PolIBCD+ and AGV^{+ins} famA DNAPs, 8-aa insertion was found solely in the famA
181 DNAP homolog of Vibrio phage ICP2 placed at the basal position of the clade of
182 PolIBCD+ and AGV^{+ins} famA DNAPs (Fig. 2). In the analyses of the second
183 alignment, AGV^{+ins} famA DNAPs grouped together with an MLBP of 92% and a BPP
184 of 0.99, excluding PolIBCD+ sequences that formed a clade with an MLBP of 72%
185 and a BPP of 0.66 (Fig. 2). The weak statistical support for the monophyly of
186 PolIBCD+ sequences is not incongruent with their paraphyletic relationship
187 reconstructed in the global famA DNAP analysis (Fig. 1B).

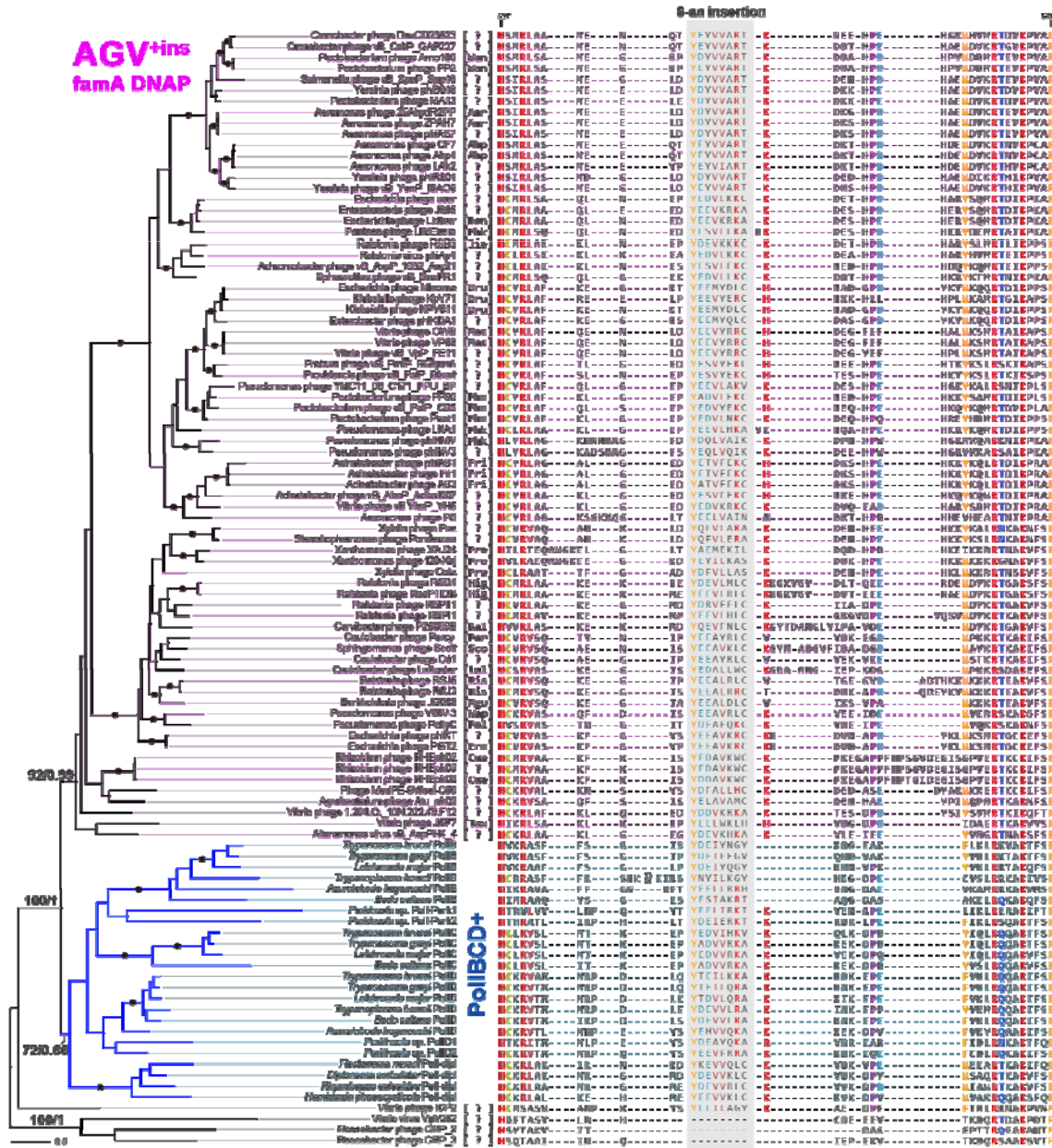


Fig. 2. Phylogenetic relationship among 74 AGV⁺ famA DNAP and 24 PoilBCD⁺ sequences that share a unique insertion of 8 amino acid residues (8-aa insertion). The tree topology and branch lengths inferred by the maximum likelihood (ML) method are shown on the left. ML bootstrap values (MLBPs) and Bayesian posterior probabilities (BPPs) for only the nodes critical to infer the origin of PoilBCD⁺ are shown. As ML and Bayesian analyses reconstructed the essentially same tree topology, only BPPs for the selected nodes are presented. The nodes supported by an MLBP of 100% and a BPP of 1.0 are marked by dots. The genus names of the autographiviruses (and podoviruses), from which famA DNAPs were sampled, are given in brackets. Abbreviations are follows: Aer, Aerosvirus; Ahp, Ahphunavirus; Bon, Bonnellvirus; Cue, Cuernavacavirus; Dru, Drulivirus; Erm, Ermolevavirus; Fri, Friunavirus; Hig, Higashivirus; Jia, Jiaoyazivirus; Kal, Kalpathivirus; Lul, Lullwatervirus; Mac, Maculvirus; Mgu, Mguuevirus; Nap, Napahavirus; Per, Percyivirus; Phk, Phikmvivirus; Phm, Phimunavirus; Pol, Pollyceevirus; Pra, Pradovirus; Ris, Risjevirus; Sco, Scottvirus; Tan, Tawavirus; Wan, Wanjuvirus; ?, unclassified. The amino acid sequences of 8-aa insertions and their flanking regions are shown on the right. 8-aa insertions are shaded in grey. The residues are colored according to their degrees of conservation. The amino acid residue numbers shown on the left and right edges of the

205 alignment are based on the famA DNAPs of Cronobacter phage DevCD23823
206 (YP_009223394.1).

207 Discussion

208 The phylogenetic analyses of the global alignment of famA DNAPs (Figs. 1A
209 and 1B) and the second alignment rich in AGV^{+ins} famA DNAP homologs (Fig. 2)
210 consistently recovered the specific affinity between PolIBCD+ and AGV^{+ins} famA
211 DNAPs. These results strongly suggest that PolIBCD+ in the extant kinetoplastids
212 and diplomonads can be traced back to a single autographivirus famA DNAP,
213 particularly the one with 8-aa insertion. In other words, PolIBCD+ is a typical
214 example of non- α -proteobacterial mt proteins established via lateral gene transfer.
215 Unfortunately, even the analyses of the second alignment, wherein the known
216 diversity of AGV^{+ins} famA DNAPs was covered, failed to pinpoint the exact origin of
217 PolIBCD+ (Fig. 2). We might be able to find an AGV^{+ins} famA DNAP homolog that
218 branches PolIBCD+ sequences directly in a future phylogenetic study covering the
219 true diversity of phage famA DNAPs. In particular, we regard that autographivirus
220 famA DNAP genes in bacterial genomes are significant. To our knowledge, no
221 autographivirus has been reported to infect eukaryotes. Thus, the common ancestor
222 of kinetoplastids and diplomonads may have acquired the famA DNAP gene from a
223 lysogenic autographivirus in a bacterial genome. If so, the bacterial genomes
224 harboring AGV^{+ins} famA DNAP genes are critical to investigate the origin of
225 PolIBCD+ at a finer level than that in the current study.

226 Members of classes Kinetoplastea and Diplomonada, together with Euglenida,
227 share another type of mitochondrion-localized famA DNAP, namely PolIA (Harada et
228 al. 2020). It is reasonable to postulate that the common ancestor of the three
229 classes—most likely the ancestral euglenozoan—had established the ancestral
230 PolIA. Although the origin of PolIA has not been addressed explicitly, past studies
231 recovered the phylogenetic link between PolIA and Pol θ , a type of famA DNAP
232 operated in the cytosol of eukaryotic cells. The original study reporting PolIA, B, C,
233 and D in *Trypanosoma brucei* has hinted at the phylogenetic affinity between PolIA
234 and Pol θ (Klingbeil et al. 2002). A recent phylogeny including famA DNAPs sampled
235 from eukaryotes and limited bacteria (Note that no phage homolog was included)
236 reconstructed a clade of PolIA and Pol θ sequences with high statistical support
237 (Harada et al. 2020). The PolIA-Pol θ affinity persisted even after the sampling of

238 famA DNAPs from bacteria and phages was improved drastically in this study (Fig.
239 1A). We here propose that the ancestral PollA was likely derived from a Pol θ
240 homolog followed by the change in subcellular localization from the cytosol to the
241 mitochondrion. Noteworthy, the evolutionary processes yielded PollA and PolIBCD+,
242 both of which are mt proteins of non- α -proteobacterial origin, are different
243 substantially from each other. The former emerged through the recycling of a pre-
244 existing eukaryotic protein while the latter is of phage origin (See above). The Pol θ
245 origin of PollA is the best estimate from both past and current phylogenetic analyses
246 of famA DNAPs but alternative possibilities still need to be explored in future studies.

247 The repertoires of mitochondrion-localized DNAPs in euglenozoans appeared
248 to be more complex than those in the majority of other eukaryotes in which a single
249 type of mitochondrion-localized DNAP (i.e. POP or Poly) seemingly operates. The
250 complexity in the repertoire of DNAPs in euglenozoan mitochondria seems to
251 coincide with that in the structure of their mtDNAs (Lukeš et al. 2002; Roy et al.
252 2007; Spencer and Gray 2011; Dobáková et al. 2015; Yabuki et al. 2016; Burger and
253 Valach 2018). Nevertheless, it is unlikely that the non- α -proteobacterial background
254 is restricted to PollA and PolIBCD+ among the proteins involved in mtDNA
255 maintenance. Rather, the machinery for mtDNA maintenance in the common
256 ancestor of kinetoplastids and diplomonids (and its descendants) are heavily
257 remodeled by both incorporating exogenous proteins via lateral gene transfer and
258 recycling the pre-existed nucleus-encoded proteins. The above conjecture can be
259 examined only after we identify the major proteins involved in DNA maintenance in
260 kinetoplastid/diplomonid mitochondria and their evolutionary origins.

261 Materials & Methods

262 Global phylogeny of famA DNAPs

263 We searched for the amino acid (aa) sequences of bacterial and phage famA
264 DNAPs in the NCBI nr database as of March 6, 2020, by BLASTP using the
265 polymerase domain of *Escherichia coli* Poll (KHH06131.1; the portion corresponding
266 to the 491st–928th aa residues) as a query (Camacho et al. 2009; Sayers et al. 2020).
267 We retrieved the sequences matched to the query with *E* values equal to or less than
268 1×10^{-4} and covered more than 200 aa in the polymerase domain. Note that the

269 sequences derived from metagenome analyses were excluded from this study. The
270 redundancy within famA DNAP sequences was removed by cluster analysis using
271 CD-HIT v4.7 with a threshold of 40% (Li and Godzik 2006; Fu et al. 2012). We finally
272 selected 119 and 327 aa sequences of phage and bacterial famA DNAPs,
273 respectively, for the downstream analyses (see below).

274 The bacterial and phage famA DNAP aa sequences (446 in total) were
275 aligned with those in eukaryotes (27 in total), namely (i) mitochondrion-localized
276 famA DNAPs in Kinetoplastea and Diplonemea (PollA, B, C, D, and Poll-dipl), (ii)
277 mitochondrion-localized famA DNAPs in animals and fungi (Poly), (iii) Pol θ localized
278 in the cytosol, (iv) mitochondrion and/or plastid-localized famA DNAPs in diverse
279 eukaryotes (POP), and (v) plastid-localized famA DNAPs in apicomplexan parasites
280 and their relatives (PREX). The aa sequences were aligned by MAFFT v7.455 with
281 the L-INS-i model (Kato and Standley 2013). Ambiguously aligned positions were
282 discarded manually, and gap-containing positions were trimmed by using trimAl v1.4
283 with the -gt 0.95 option (Capella-Gutiérrez et al. 2009). The final "global famA DNAP"
284 alignment comprised 473 sequences with 316 unambiguously aligned aa positions.
285 The final global famA alignment is provided as a part of the supplementary materials.
286 We subjected this alignment to the ML phylogenetic analysis by IQ-TREE v1.6.12
287 using the LG + Γ + F + C60 + PMSF model (Nguyen et al. 2015; Wang et al. 2018).
288 The guide tree was obtained using the LG + Γ + F model that was selected by
289 ModelFinder (Kalyaanamoorthy et al. 2017). The statistical support for each
290 bipartition in the ML tree was calculated by 100-replicate non-parametric bootstrap
291 analysis.

292 Phylogenetic analyses of an alignment rich in autographivirus famA 293 DNAPs

294 We retrieved 175 famA DNAP aa sequences of autographiviruses from the NCBI nr
295 database. The details of the survey were the same as described above. The 175
296 famA DNAPs were sampled from 99 members belonging to 57 genera and 76
297 unclassified members in the family Autographiviridae. The autographivirus famA
298 DNAPs were found to comprise two types based on the presence/absence of an
299 insertion of 8 aa residues (8-aa insertion; see above). The famA DNAPs with 8-aa
300 insertion (AGV^{+ins} famA DNAPs) appeared to be closely related to PollBCD+,

301 mitochondrion-localized famA DNAPs in kinetoplastids (PollB, C, D, Poll-Perk1/2)
302 and that in diplomonads (Poll-dipl). The redundancy among the AGV^{+ins} famA DNAPs
303 was reduced by cluster analysis using CD-HIT v4.7 with a threshold of 90%. Finally,
304 we aligned the aa sequences of 74 AGV^{+ins} famA DNAPs, 24 PollBCD+, and famA
305 DNAPs of four members of Podoviridae by MAFFT v7.455 with the L-INS-i model.
306 Ambiguously aligned positions were discarded manually, and gap-containing
307 positions were trimmed by using trimAl v1.4 with the -gt 0.9 option. The final version
308 of the second alignment is provided as a part of the supplementary materials. The
309 final alignment containing 102 sequences with 581 unambiguously aligned aa
310 positions was subjected to both ML and Bayesian phylogenetic analyses. The ML
311 and ML bootstrap analyses were performed as described above. For Bayesian
312 analysis using Phylobayes v4.1, we run four Markov Chain Monte Carlo chains for
313 100,000 cycles with burn-in of 25,000 (maxdiff = 0.09472) and calculated the
314 consensus tree with branch lengths and BPPs from the remaining trees (Lartillot et al.
315 2009). The amino acid substitution model was set to CAT + GTR in Phylobayes
316 analysis described above.

317 Data availability

318 The alignment datasets for phylogenetic analysis are available in supplementary
319 materials at
320 [https://drive.google.com/drive/folders/1vpwh0MzYul_wjKmyutZIZR1MSmMZn5ca?us](https://drive.google.com/drive/folders/1vpwh0MzYul_wjKmyutZIZR1MSmMZn5ca?usp=sharing)
321 [p=sharing](https://drive.google.com/drive/folders/1vpwh0MzYul_wjKmyutZIZR1MSmMZn5ca?usp=sharing).

322 Acknowledgments

323 This work was supported by the grants from the Japanese Society for Promotion of
324 Sciences awarded to Y. I. (numbers 18KK0203 and 19H03280).

325 Literature cited

326 Adriaenssens EM, Sullivan MB, Knezevic P, van Zyl LJ, Sarkar BL, et al. 2020.
327 Taxonomy of prokaryotic viruses: 2018-2019 update from the ICTV
328 bacterial and archaeal viruses subcommittee. *Arch Virol.* 165(5):1253–
329 1260. doi:10.1007/s00705-020-04577-8.

- 330 Burger G, Valach M. 2018. Perfection of eccentricity: mitochondrial genomes of
331 diplonemids. *IUBMB Life*. 70(12):1197–1206. doi:10.1002/iub.1927.
- 332 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, et al. 2009. BLAST+:
333 Architecture and applications. *BMC Bioinformatics*. 10:421:1-421:9.
334 doi:10.1186/1471-2105-10-421.
- 335 Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. 2009. trimAl: a tool for
336 automated alignment trimming in large-scale phylogenetic analyses.
337 *Bioinformatics*. 25(15):1972–1973. doi:10.1093/bioinformatics/btp348.
- 338 Dobáková E, Flegontov P, Skalický T, Lukeš J. 2015. Unexpectedly streamlined
339 mitochondrial genome of the euglenozoan *Euglena gracilis*. *Genome Biol*
340 *Evol*. 7(12):3358–3367. doi:10.1093/gbe/evv229.
- 341 El-Gebali S, Mistry J, Bateman A, Eddy SR, Luciani A, et al. 2019. The Pfam protein
342 families database in 2019. *Nucleic Acids Res*. 47(D1):D427–D432.
343 doi:10.1093/nar/gky995.
- 344 Fu L, Niu B, Zhu Z, Wu S, Li W. 2012. CD-HIT: accelerated for clustering the next-
345 generation sequencing data. *Bioinformatics*. 28(23):3150–3152.
346 doi:10.1093/bioinformatics/bts565.
- 347 Gabaldón T, Huynen MA. 2007. From endosymbiont to host-controlled organelle: the
348 hijacking of mitochondrial protein synthesis and metabolism. *PLoS*
349 *Comput Biol*. 3(11):2209–2218. doi:10.1371/journal.pcbi.0030219.
- 350 Gray MW. 2015. Mosaic nature of the mitochondrial proteome: implications for the
351 origin and evolution of mitochondria. *Proc Natl Acad Sci U S A*.
352 112(33):10133–10138. doi:10.1073/pnas.1421379112.
- 353 Graziewicz MA, Longley MJ, Copeland WC. 2006. DNA polymerase γ in
354 mitochondrial DNA replication and repair. *Chem Rev*. 106(2):383–405.
355 doi:10.1021/cr040463d.
- 356 Harada R, Hirakawa Y, Yabuki A, Kashiya Y, Maruyama M, et al. 2020. Inventory
357 and evolution of mitochondrion-localized family A DNA polymerases in
358 Euglenozoa. *Pathogens*. 9(4):257. doi:10.3390/pathogens9040257.
- 359 Hirakawa Y, Watanabe A. 2019. Organellar DNA polymerases in complex plastid-
360 bearing algae. *Biomolecules*. 9(4):140:1-140:12.
361 doi:10.3390/biom9040140.
- 362 Janouškovec J, Tikhonenkov D V., Burki F, Howe AT, Kolísko M, et al. 2015. Factors
363 mediating plastid dependency and the origins of parasitism in

- 364 apicomplexans and their close relatives. *Proc Natl Acad Sci U S A*.
365 112(33):10200–10207. doi:10.1073/pnas.1423790112.
- 366 Jones P, Binns D, Chang HY, Fraser M, Li W, et al. 2014. InterProScan 5: genome-
367 scale protein function classification. *Bioinformatics*. 30(9):1236–1240.
368 doi:10.1093/bioinformatics/btu031.
- 369 Jung GH, Leavitt MC, Hsieh JC, Ito J. 1987. Bacteriophage PRD1 DNA polymerase:
370 evolution of DNA polymerases. *Proc Natl Acad Sci U S A*. 84(23):8287–
371 8291. doi:10.1073/pnas.84.23.8287.
- 372 Kalyaanamoorthy S, Minh BQ, Wong TKF, Von Haeseler A, Jermini LS. 2017.
373 ModelFinder: fast model selection for accurate phylogenetic estimates.
374 *Nat Methods*. 14(6):587–589. doi:10.1038/nmeth.4285.
- 375 Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version
376 7: improvements in performance and usability. *Mol Biol Evol*. 30(4):772–
377 780. doi:10.1093/molbev/mst010.
- 378 Klingbeil MM, Motyka SA, Englund PT. 2002. Multiple mitochondrial DNA
379 polymerases in *Trypanosoma brucei*. *Mol Cell*. 10(1):175–186.
380 doi:10.1016/S1097-2765(02)00571-3.
- 381 Krasich R, Copeland WC. 2017. DNA polymerases in the mitochondria: a critical
382 review of the evidence. *Physiol Behav*. 22(1):692–709.
383 doi:10.1016/j.physbeh.2017.03.040.
- 384 Lartillot N, Lepage T, Blanquart S. 2009. PhyloBayes 3: a bayesian software
385 package for phylogenetic reconstruction and molecular dating.
386 *Bioinformatics*. 25(17):2286–2288. doi:10.1093/bioinformatics/btp368.
- 387 Lavigne R, Seto D, Mahadevan P, Ackermann HW, Kropinski AM. 2008. Unifying
388 classical and molecular taxonomic classification: analysis of the
389 Podoviridae using BLASTP-based tools. *Res Microbiol*. 159(5):406–414.
390 doi:10.1016/j.resmic.2008.03.005.
- 391 Li W, Godzik A. 2006. Cd-hit: a fast program for clustering and comparing large sets
392 of protein or nucleotide sequences. *Bioinformatics*. 22(13):1658–1659.
393 doi:10.1093/bioinformatics/btl158.
- 394 Lukeš J, Guilbride DL, Votýpka J, Zíková A, Benne R, et al. 2002. Kinetoplast DNA
395 network: evolution of an improbable structure. *Eukaryotic Cell*.
396 1(4):495–502. doi:10.1128/EC.1.4.495.

- 397 Moriyama T, Terasawa K, Sato N. 2011. Conservation of POPs, the plant organellar
398 DNA polymerases, in eukaryotes. *Protist.* 162(1):177–187.
399 doi:10.1016/j.protis.2010.06.001.
400 <http://dx.doi.org/10.1016/j.protis.2010.06.001>.
- 401 Nguyen LT, Schmidt HA, Von Haeseler A, Minh BQ. 2015. IQ-TREE: a fast and
402 effective stochastic algorithm for estimating maximum-likelihood
403 phylogenies. *Mol Biol Evol.* 32(1):268–274. doi:10.1093/molbev/msu300.
- 404 Roger AJ, Muñoz-Gómez SA, Kamikawa R. 2017. The origin and diversification of
405 mitochondria. *Curr Biol.* 27(21):R1177–R1192.
406 doi:10.1016/j.cub.2017.09.015.
- 407 Roy J, Faktorová D, Lukeš J, Burger G. 2007. Unusual mitochondrial genome
408 structures throughout the Euglenozoa. *Protist.* 158(3):385–396.
409 doi:10.1016/j.protis.2007.03.002.
- 410 Sayers EW, Beck J, Brister JR, Bolton EE, Canese K, et al. 2020. Database
411 resources of the national center for biotechnology information. *Nucleic*
412 *Acids Res.* 48(D1):D9–D16. doi:10.1093/nar/gkz899.
- 413 Spencer DF, Gray MW. 2011. Ribosomal RNA genes in *Euglena gracilis*
414 mitochondrial DNA: fragmented genes in a seemingly fragmented
415 genome. *Mol Genet Genomics.* 285(1):19–31. doi:10.1007/s00438-010-
416 0585-9.
- 417 Wang HC, Minh BQ, Susko E, Roger AJ. 2018. Modeling site heterogeneity with
418 posterior mean site frequency profiles accelerates accurate
419 phylogenomic estimation. *Syst Biol.* 67(2):216–235.
420 doi:10.1093/sysbio/syx068.
- 421 Wang Z, Wu M. 2014. Phylogenomic reconstruction indicates mitochondrial ancestor
422 was an energy parasite. *PLoS One.* 9(10):e110685.
423 doi:10.1371/journal.pone.0110685.
- 424 Yabuki A, Tanifuji G, Kusaka C, Takishita K, Fujikura K. 2016. Hyper-Eccentric
425 structural genes in the mitochondrial genome of the algal parasite
426 *Hemistasia phaeocysticola*. *Genome Biol Evol.* 8(9):2870–2878.
427 doi:10.1093/gbe/evw207.