

1 **Title**

2 The economics of endosymbiotic gene transfer and the evolution of organellar genomes

3 **One sentence summary**

4 The high copy number of organellar genomes renders endosymbiotic gene transfer energetically
5 favourable for the vast majority of organellar genes.

6 **Authors**

7 Steven Kelly

8 **Affiliations**

9 1) Department of Plant Sciences, University of Oxford, South Parks Road, OX1 3RB

10 **ORCID ID**

11 SK: 0000-0001-8583-5362

12 **Corresponding author**

13 Name: Steven Kelly

14 Email: steven.kelly@plants.ox.ac.uk

15 Telephone: +44 (0)1865 275123

16 **Abstract**

17 The endosymbiosis of the bacterial progenitors of mitochondrion and the chloroplast are landmark
18 events in the evolution of life on earth. While both organelles have retained substantial proteomic
19 and biochemical complexity, this complexity is not reflected in the content of their genomes.
20 Instead, the organellar genomes encode fewer than 5% of genes found in close relatives of their
21 ancestors. While some of the 95% of missing organellar genes have been discarded, many have
22 been transferred to the host nuclear genome through a process known as endosymbiotic gene
23 transfer. Here we demonstrate that the energy liberated or consumed by a cell as a result of
24 endosymbiotic gene transfer is sufficient to provide a selectable advantage for retention or nuclear-
25 transfer of organellar genes in eukaryotic cells. We further demonstrate that for realistic estimates
26 of protein abundances, organellar protein import costs, host cell sizes, and cellular investment in
27 organelles that it is energetically favourable to transfer the majority of organellar genes to the
28 nuclear genome. Moreover, we show that the selective advantage of such transfers is sufficiently
29 large to enable such events to rapidly reach fixation. Thus, endosymbiotic gene transfer can be

30 advantageous in the absence of any additional benefit to the host cell, providing new insight into
31 the processes that have shaped eukaryotic genome evolution.

32 **Main**

33 Endosymbiosis has underpinned two of the most important innovations in the history of life on
34 Earth (Archibald 2015a; Martin, et al. 2015). The endosymbiosis of the alphaproteobacterium that
35 became the mitochondrion led to the emergence and radiation of the eukaryotes (Yang, et al.
36 1985; Martin and Müller 1998; Roger, et al. 2017), and the endosymbiosis of the cyanobacterium
37 that became the chloroplast first enabled oxygenic photosynthesis in eukaryotes (Martin and
38 Kowallik 1999; Archibald 2015b). The function and evolution of both organelles is inextricably
39 linked with energy metabolism and the evolution of the eukaryotic cell (Lane and Martin 2010; Lane
40 2014; Booth and Doolittle 2015a, b; Lane and Martin 2015; Lynch and Marinov 2017; Roger, et al.
41 2017; Lynch and Marinov 2018), and has given rise to the multicellular organisms that dominate
42 the biosphere (Bar-On, et al. 2018). Following both of these endosymbioses there was a dramatic
43 reduction in the gene content of the endosymbiont genomes such that extant mitochondria and
44 chloroplasts typically harbour fewer than 5% of the genes found in their free-living prokaryotic
45 relatives (Gray, et al. 1999; Timmis, et al. 2004; Green 2011). While many of the original
46 endosymbiont genes have been lost through mutation and drift (Lynch, et al. 2006; McCutcheon
47 and Moran 2012; Smith and Keeling 2015; Smith 2016), others have been transferred to the host
48 nuclear genome and their products imported back into the organelle where they function (Martin, et
49 al. 2002; Brown 2003; Deusch, et al. 2008; Thiergart, et al. 2012; Dagan, et al. 2013). For
50 example, the mitochondrion of humans (Calvo and Mootha 2010) and chloroplasts of plants (Ferro,
51 et al. 2010) each contain more than 1000 proteins yet their genomes encode fewer than 100
52 genes. Therefore, the reduced gene content of organelles is not representative of their molecular,
53 proteomic or biochemical complexity. Furthermore, endosymbiotic gene transfer is not unique to
54 the evolution of chloroplasts and mitochondria but has also been observed with bacterial
55 endosymbionts of insects (McCutcheon and Moran 2012; Husnik, et al. 2013) and with the
56 endosymbiosis of the chromatophore of *Paulinella* (Nakayama and Ishida 2009; Nowack, et al.
57 2010; Reyes-Prieto, et al. 2010; Singer, et al. 2017; Nowack and Weber 2018). Thus,

58 endosymbiont genome reduction and endosymbiotic gene transfer are recurring themes in the
59 evolution of eukaryotic nuclear and cytoplasmic genomes.

60 Given, its fundamental importance to the evolution of eukaryotic genomes, several hypotheses
61 have been proposed to explain why endosymbiotic gene transfer occurs (Herrmann 1997; Martin
62 and Herrmann 1998; Daley and Whelan 2005; Reyes-Prieto, et al. 2006; Speijer, et al. 2020). For
63 example, it has been proposed that it protects endosymbiont genes from mutational hazard (Allen
64 and Raven 1996; Lynch, et al. 2006; Smith 2016; Speijer, et al. 2020), and that it enables
65 endosymbiont genes that are otherwise trapped in a haploid genome to recombine and thus
66 escape from Muller's ratchet (Muller 1964; Lynch 1996; Martin and Herrmann 1998; Lynch, et al.
67 2006; Neiman and Taylor 2009; Smith 2016). It has also been proposed that endosymbiotic gene
68 transfer is an inevitable consequence of a constant stream of endosymbiont genes entering the
69 nucleus (Doolittle 1998), and that transfer to the nuclear genome allows the host cell to gain better
70 control over the replication and function of the organelle (Herrmann 1997) allowing wider cellular
71 network integration (Nowack, et al. 2010; Reyes-Prieto 2015). However, mutation rates of
72 organellar genes are often not higher than nuclear genes (Wolfe, et al. 1987; Lynch, et al. 2006;
73 Lynch, et al. 2007; Drouin, et al. 2008; Smith 2015; Smith and Keeling 2015; Smith 2016; Grisdale,
74 et al. 2019) and therefore effective mechanisms for protection against DNA damage in organelles
75 must exist. Similarly, although there is evidence for the action of Muller's ratchet in mitochondria
76 (Lynch 1996; Neiman and Taylor 2009) chloroplasts appear largely to escape this effect (Wolfe, et
77 al. 1987; Lynch 1997) likely due to gene conversion (Khakhlova and Bock 2006), and thus it does
78 not fully explain why endosymbiotic gene transfer occurred in both lineages. Finally, the nature of
79 the regulatory advantage for having genes reside in the nuclear genome is difficult to quantify, and
80 may simply be a projection of anthropocentric ideals of centralised control onto the nucleus of the
81 host cell. Thus, it is unclear whether endosymbiotic gene transfer functions simply as rescue from
82 processes that would otherwise lead to gene loss, or whether there may also be an advantage to
83 the cell for retaining an endosymbiont gene to the nuclear genome.

84 We hypothesised that an advantage for endosymbiotic gene transfer may arise from the difference
85 in the cost to the cell of encoding a gene in the organellar and nuclear genome. This is because

86 each eukaryotic cell typically contains multiple organelles and each organelle typically harbours
87 multiple copies of the organellar genome (Bendich 1987; Cole 2016). The number of organelles in
88 a cell reflects the biochemical requirement of that cell for those organelles, and the high genome
89 copy number per organelle has been proposed to provide protection against DNA damage
90 (Shokolenko, et al. 2009) and to enable the organelle to achieve high protein abundance for genes
91 encoded in the organellar genome (Bendich 1987). Thus, while a typical diploid eukaryotic cell
92 contains two copies of the nuclear genome, the same cell contains hundreds to hundreds of
93 thousands of copies of its organellar genomes (Bendich 1987; Cole 2016). As DNA costs energy
94 and cellular resources to biosynthesise (Lynch and Marinov 2015), the cost to the cell of encoding
95 a gene in the organellar and nuclear genome is different. To quantify this difference, we evaluated
96 the cost of encoding a gene in the nuclear or organellar genome. Here, the cost of a gene was
97 considered to be the cost of the chromosome divided by the number of genes on that chromosome
98 to account for introns, structural, and regulatory elements (we also included the cost of the
99 requisite number of histone proteins contained in nucleosomes for nuclear genes). This revealed
100 that the cost of encoding a gene in the organellar genome is on average one order of magnitude
101 higher than the cost of encoding a gene in the nuclear genome (Figure 1A). This difference is
102 further enhanced if the biosynthesis cost of just the coding sequences of the genes are compared
103 directly (Figure 1B). Thus, the cost to the cell of encoding a gene in the organellar genome is
104 substantially higher than the cost of encoding the same gene in the nuclear genome.
105 Consequently, for any essential organellar gene the cell may be able to save resources by
106 transferring that gene from the organellar genome to the nuclear genome. For example,
107 endosymbiotic transfer of a 1000 bp gene from the mitochondrion to the nuclear genome in
108 humans, yeast or *Arabidopsis* would save 5,000,000 bp, 200,000 bp or 100,000 bp of DNA per
109 cell, respectively, and an analogous transfer from the chloroplast genome to the nuclear genome in
110 *Arabidopsis* would save 1,500,000 bp of DNA per cell. We hypothesised that if the energy saved
111 by transferring such a gene offset the cost of importing the required abundance of gene product
112 back into the organelle then this would provide a direct energetic and fitness advantage to the host
113 cell for endosymbiotic gene transfer.

114 To test this hypothesis, we assessed the conditions under which it is more energetically favourable
115 to encode a gene in the organellar or nuclear genome. Here, the free energy of endosymbiotic
116 gene transfer (which we define as the difference in energy cost between a cell which encodes a
117 given gene in the organellar genome and a cell which encodes the same gene in the nuclear
118 genome and imports the requisite amount of gene product into to the organelle, see Methods) was
119 computed for an average length bacterial gene as a function of protein abundance, protein import
120 cost, and organellar genome copy number. This revealed that there is a simple relationship such
121 that the higher the copy number of the organellar genome, the more energy that is liberated by
122 endosymbiotic gene transfer and thus the more protein that can be imported into the organelle
123 while still reducing the overall energetic cost of the cell (Figure 2A). To simulate the organellar
124 genome reduction that would result if all such energetically favourable endosymbiotic gene
125 transfers occurred, the complete genomes with measured protein abundances for an
126 alphaproteobacterium (*Bartonella henselae*) and a cyanobacterium (*Microcystis aeruginosa*) were
127 subject to a simulated endosymbiosis. Here, a range of host cell sizes was simulated such that
128 they encompassed the majority of diversity exhibited by extant eukaryotes (Milo 2013) and would
129 thus likely encompass the size range of the host cell that originally engulfed the
130 alphaproteobacterial and cyanobacterial organellar progenitors. This range extended from a small
131 unicellular yeast-like cell (10^7 proteins), to a typical unicellular algal cell (10^8 proteins) to a large
132 metazoan/plant cell (10^9 proteins). Each of these cell types were then considered to allocate a
133 realistic range of total cellular protein to mitochondria/chloroplasts representative of extant
134 eukaryotic cells (Supplemental Table S1). For each simulated endosymbiosis, the free energy of
135 endosymbiotic gene transfer was calculated for each gene given its measured protein abundance
136 (Wang, et al. 2015) and a realistic range of protein import costs (including the total biosynthetic
137 cost of the protein import machinery, See Methods). This revealed that for a broad range of
138 estimates of cell size, organellar genome copy number, organellar fraction (i.e. the fraction of the
139 total number of protein molecules in a cell that are contained within the organelle), protein
140 abundance, and protein import cost it is energetically favourable to the cell to transfer the majority
141 of organellar genes to the nuclear genome and re-import the proteins back to the organelle (Figure
142 2B and 2C). Here, only the proteins with the highest abundance, and thus which occur the largest

143 import cost, are retained in the organellar genomes. While other examples of eukaryotic cell sizes
144 and resource allocation outside the range shown here exist in nature, and the properties of the cell
145 which engulfed the progenitors of the mitochondrion and chloroplast are unknown, the properties of
146 the cells are likely encompassed within the ranges presented here.

147 To estimate the strength of selection that would act on the change in energy incurred from an
148 endosymbiotic gene transfer event, the free energy of endosymbiotic gene transfer for each gene
149 was placed in context of the total energy budget of the host cell. As above, this analysis was
150 conducted for a broad range of host cell size, organellar fraction, endosymbiont genome copy
151 number, and protein import cost that is representative of a broad range of eukaryotic cells (Figure
152 3A and B, Supplemental Figures S1 – S6, Supplemental Table S2). This revealed that for even
153 modest per-cell endosymbiont genome copy numbers (≥ 100 copies per cell) the selection
154 coefficients for the transfer of the majority of endosymbiont genes are relatively large $\sim 1 \times 10^{-4}$
155 (Figure 3, Supplemental Figures S1 – S6), $\sim 10,000$ times stronger than the selection coefficient
156 acting against disfavoured synonymous codons (Hartl, et al. 1994). Moreover, for high per-cell
157 endosymbiont genome copy numbers (≥ 1000 genome copies per cell) these selection coefficients
158 are large ($\sim 1 \times 10^{-3}$) and similar to the strength of selection that caused the allele conferring lactose
159 tolerance to rapidly sweep through human populations in ~ 500 generations (Bersaglieri, et al.
160 2004). In contrast, selection coefficients for retention of genes in the organellar genome only occur
161 when organellar genome copy numbers are low, and/or when large proportions of cellular
162 resources are invested in organelle (Figure 3A and B, Supplemental Figures S1 – S6). However,
163 with the exception of very highly abundant proteins (discussed below) these selection coefficients
164 are generally weaker. Thus, over a broad range of host cell sizes, organellar genome copy
165 numbers, organellar fractions, and per-protein ATP import costs, endosymbiotic gene transfer of
166 the majority of genes is sufficiently energetically advantageous that any such transfer events, if
167 they occurred, would rapidly reach fixation (Supplemental Figure S7). Thus, endosymbiotic gene
168 transfer is intrinsically advantageous to the cell for the majority of organellar genes in the absence
169 of additional benefits.

170 Although the free energy of endosymbiotic gene transfer is sufficient to explain why organellar
171 genes are transferred to the nucleus, it is not proposed that it is the only factor that influences the
172 location of an organellar gene. Instead, a large cohort of factors including the requirement for
173 organellar mediated RNA editing, protein chaperones, protein folding, post-translational
174 modifications, escaping mutation hazard, Muller's ratchet, enhanced nuclear control, and drift will
175 act antagonistically or synergistically with the free energy of endosymbiotic gene transfer to
176 influence the set of genes that are retained in, or transferred from, the organellar genomes.
177 Moreover, the free energy of endosymbiotic gene transfer provides a mechanistic basis for
178 selection to act for or against Doolittle's "You are what you eat" ratchet for endosymbiotic gene
179 transfer (Doolittle 1998). It is noteworthy in these contexts, that if the protein encoded by the
180 endosymbiont gene can provide its function outside of the endosymbiont (e.g. by catalysing a
181 reaction that could occur equally well in the cytosol of the host as in the endosymbiont) then the
182 energetic advantage of gene transfer to the nuclear genome is further enhanced, as the cost of
183 protein import is not incurred. Similarly, although gene loss is predominantly thought to be
184 mediated by mutation pressure and drift (Lynch, et al. 2006), the elevated per-cell endosymbiont
185 genome copy number also provides an energetic incentive to the host cell for complete gene loss.
186 Thus, the high genome copy number required to protect DNA from damage (Shokolenko, et al.
187 2009) and facilitate high levels of protein production (Bendich 1987), also provides the energetic
188 incentive for the cell to delete endosymbiont genes as well as transfer them to the nuclear genome.
189 The analysis presented here shows that for a broad range of cell sizes and resource allocations
190 that endosymbiotic gene transfer of the majority of organellar genes is energetically favourable and
191 thus advantageous to the cell. Retention of genes in the organellar genomes is only favourable
192 under conditions where the encoded organellar protein is required in very high abundance and/or
193 the copy number of the organellar genome is low (Figure 2B, 2C, 3A and 3B). The interaction
194 between protein abundance and genome copy number provides some insight into why organellar
195 genomes still retain some genes. For example, in large plant cells such as those in the leaves of
196 *Arabidopsis thaliana* it is unfavourable to transfer the *rbcL* gene encoding the RuBisCO large
197 subunit from the chloroplast genome to the nuclear genome, as although it would save 8.7×10^7
198 ATP per cell in DNA biosynthesis costs it would incur a daily cost of $\sim 3.96 \times 10^{12}$ ATP per cell

199 (0.17% of the daily energy budget of the cell) just to import the required amount of RuBisCO large
200 subunit back into the chloroplast (see methods). Thus, from a cost perspective it is energetically
201 favourable to retain this gene in the chloroplast genome. The same is also true for 62 of the 88
202 genes currently found in the chloroplast genome in *Arabidopsis thaliana* (Supplemental Table S3)
203 such that selection would act against transfer of these genes from the chloroplast genome. In,
204 contrast it is energetically favourable to transfer the majority of genes from the mitochondrial
205 genome to the nuclear genome in *Arabidopsis* (99 out of 122), and all of the genes encoded in the
206 human mitochondrial genome to the human nuclear genome (Supplemental Table S3). Thus, high
207 cellular investment in chloroplast proteins creates a selectable advantage for retention of the
208 majority of genes currently encoded in the chloroplast genome.

209 While we do not know precisely what the cells that engulfed the progenitors of the mitochondrion or
210 the chloroplast looked like (as only extant derivatives survive), it is safe to assume that cell size
211 and investment in organelles has altered since these primary endosymbioses first occurred.
212 Accordingly, the selective advantage (or disadvantage) of transfer of any given gene is transient
213 and will have varied during the radiation of the eukaryotes as cell size and organellar volume
214 evolved and changed in disparate eukaryotic lineages. This coupled with the lack of an organellar
215 protein export system (i.e. from the organelle to the host cytosol) and the presence (and
216 acquisition) of introns in nuclear encoded genes (Rogozin, et al. 2012) means that it is more
217 difficult for endosymbiotic gene transfer to operate in the reverse direction (i.e. from the nucleus to
218 organelle). Collectively, this would create a ratchet-like effect trapping genes in the nuclear
219 genome even if subsequent changes in cell size and investment in organelles means that it
220 became energetically advantageous to return the gene to the organelle later in evolution. Thus,
221 current organellar and nuclear gene contents predominantly reflect past pressures to transfer
222 genes to the nuclear genome.

223 Endosymbiotic gene transfer is a recurring theme in the evolution of the eukaryotic tree of life. The
224 discovery that the free energy of endosymbiotic gene transfer can act to promote retention or
225 transfer of organellar genes to the nuclear genome uncovers a novel process that has helped
226 shape the content and evolution of organellar and nuclear genomes in eukaryotes. Moreover, it

227 helps to explain why organelles have surrendered the vast majority of their genes for the sake of
228 the greater good of the cell.

229 **Materials and Methods**

230 **Data sources**

231 The *Arabidopsis thaliana* genome sequence and corresponding set of representative gene models
232 were downloaded from Phytozome V13 (Goodstein, et al. 2012). The human genome sequence
233 and gene models from assembly version GRCh38.p13 (GCA_000001405.28), the *Bartonella*
234 *henselae* genome sequence and gene models from assembly version ASM4670v1, the *Microcystis*
235 *aeruginosa* NIES-843 genome sequence and gene models from assembly version ASM1062v1
236 were each downloaded from Ensembl (Yates, et al. 2020). The *Saccharomyces cerevisiae*
237 sequence and gene models from assembly version R64-2-1_20150113 were downloaded from the
238 *Saccharomyces* Genome Database (Cherry, et al. 2012). Protein abundance data for all species
239 were obtained from PAXdb v4.1 (Wang, et al. 2015).

240 **Constants used to evaluate the per cell ATP costs of genes and chromosomes**

241 The ATP biosynthesis cost of nucleotides and amino acids was obtained from (Chen, et al. 2016)
242 and (Lynch and Marinov 2015) and are provided in Supplemental Table S4. *The Homo sapiens*
243 mitochondrial genome copy number of 5000 was obtained from (Cole 2016). The *Saccharomyces*
244 *cerevisiae* mitochondrial genome copy number of 200 was obtained from (Miyakawa 2017). The
245 *Arabidopsis thaliana* chloroplast genome copy number of 1500 was obtained from (Zoschke, et al.
246 2007) and the *Arabidopsis thaliana* mitochondrial genome copy number of 100 was obtained from
247 (Cole 2016).

248 For genes in nuclear chromosomes the cost of DNA was calculated to include the cost of
249 nucleosomes with one histone octamer comprising two copies each of the histone proteins H2A,
250 H2B, H3, and H4 every 180bp (147bp for the two turns of DNA around the histone octamer and
251 33bp for the spacer) (Lynch and Marinov 2015). For organellar chromosomes there are no
252 histones/nucleosomes and thus the biosynthetic cost of genes in organellar chromosomes was
253 calculated as cost of the DNA divided by the number of genes on the chromosome (Supplemental
254 Table S5).

255 The average gene length used for the simulation study in Figure 2 was obtained by computing the
256 average gene length across the two bacterial genomes used in this study, *Bartonella henselae*
257 ASM4670v1 and *Microcystis aeruginosa* NIES-843.

258 **Calculating protein import costs**

259 Although the molecular mechanisms of mitochondrial and chloroplast protein import differ (Soll and
260 Schleiff 2004; Jarvis 2008; Wiedemann and Pfanner 2017) they share many commonalities
261 including the requirement for energy in the form of nucleoside triphosphate hydrolysis (Schatz and
262 Dobberstein 1996). The energetic cost of mitochondrial or chloroplast protein import is difficult to
263 measure directly, and accordingly estimates vary over two orders of magnitude from ~0.05 ATP
264 per amino acid to 5 ATP per amino acid (Mokranjac and Neupert 2008; Shi and Theg 2013;
265 Backes and Herrmann 2017). Thus, for the purposes of this study the full range of estimates was
266 considered in all simulations when evaluating the import cost of organellar targeted proteins
267 encoded by nuclear genes.

268 The cost of the biosynthesis of the protein import machinery (i.e. the TOC/TIC or TOM/TIM
269 complexes, Supplemental Table S6) was also included in the per protein import costs calculated in
270 this study. For *Arabidopsis thaliana*, if the total ATP biosynthesis cost of all TOC/TIC complex
271 proteins in the cell (i.e. the full biosynthesis cost of all the amino acids of all the proteins at their
272 measured abundance in the cell) is distributed equally among all of the proteins that are imported
273 into the chloroplast then it would add an additional 0.2 ATP per residue imported (Supplemental
274 Table S7). Similarly, if the total ATP biosynthesis cost of all TOM/TIM proteins in the cell in *Homo*
275 *sapiens*, *Saccharomyces cerevisiae* and *Arabidopsis thaliana* is distributed equally among all of
276 the proteins that are imported into the mitochondrion in those species then it would add an
277 additional 0.2 ATP, 0.7 ATP, and 0.2 ATP per residue imported, respectively (Supplemental Table
278 S7). In all cases the proteins that were predicted to be imported into the organelle were identified
279 using TargetP-2.0 (Almagro Armenteros, et al. 2019) and protein abundance was calculated using
280 measured protein abundance estimates for each species obtained from PAXdb 4.0 (Wang, et al.
281 2015), assuming a total cell protein content of 1×10^9 proteins for a human cell, 1×10^7 proteins for a
282 yeast cell and 2.5×10^{10} proteins for an *Arabidopsis thaliana* cell. As we modelled ATP import

283 costs from 0.05 ATP to 5 ATP per-residue the cost of the import machinery was considered to be
284 included within the bounds considered in this analysis.

285 **Evaluating the proportion of the total proteome invested in organelles**

286 To provide estimates of the fraction of cellular protein resources invested in organellar proteomes
287 the complete predicted proteomes and corresponding protein abundances were quantified.
288 Organellar targeting was predicted using TargetP-2.0 (Almagro Armenteros, et al. 2019) and
289 protein abundance estimates obtained from PAXdb 4.0 (Wang, et al. 2015). The proportion of
290 cellular resources are provided in Supplemental Table S1 and were used to provide the indicative
291 regions or parameter space occupied by metazoa, yeast and plants shown on Figure 2B and C.
292 Specifically, ~5% of total cellular protein is contained within mitochondria in *H. sapiens*, *S.*
293 *cerevisiae* and *A. thaliana* and ~50% of total cellular protein is contained within chloroplasts in *A.*
294 *thaliana*.

295 **Calculating the free energy of endosymbiotic gene transfer**

296 The free energy of endosymbiotic gene transfer (ΔE_{EGT}) is evaluated as the difference in ATP
297 biosynthesis cost required to encode a gene (ΔD) in the endosymbiont genome (D_{end}) and the
298 nuclear genome (D_{nuc}) minus the difference in ATP biosynthesis cost required to produce the
299 protein (ΔP) in the organelle (P_{end}) vs in the cytosol (P_{cyt}) and ATP cost to import the protein into
300 the organelle (P_{import}). Such that

$$301 \quad \Delta E_{EGT} = \Delta D - \Delta P \quad [1]$$

302 Where

$$303 \quad \Delta D = D_{end} - D_{nuc} \quad [2]$$

304 And

$$305 \quad \Delta P = P_{end} - P_{cyt} - P_{import} \quad [3]$$

306 The energetic cost of producing a protein in the endosymbiont and in the cytosol are assumed to
307 be equal and thus

$$308 \quad \Delta P = P_{import} \quad [4]$$

309 P_{import} is evaluated as the product of the product of the length of the amino acid sequence (L_{prot}),
310 the ATP cost of importing a single residue from the contiguous polypeptide chain of that protein
311 (C_{import}), the number of copies of that protein contained within the cell that must be imported (N_p)
312 such that

$$313 \quad \Delta P = P_{import} = L_{prot} C_{import} N_p \text{ [5]}$$

314 Both D_{end} and D_{nuc} are evaluated as the product of the ATP biosynthesis cost of the double
315 stranded DNA (A_{DNA}) that comprises the gene under consideration and the copy number (C) of the
316 genome in the cell such that

$$317 \quad D_{end} = A_{DNA} C_{end} \text{ [6]}$$

318 And

$$319 \quad D_{nuc} = A_{DNA} C_{nuc} \text{ [7]}$$

320 Such that

$$321 \quad \Delta D = A_{DNA} (C_{end} - C_{nuc}) \text{ [8]}$$

322 Where C_{end} and C_{nuc} are the per-cell copy number of the endosymbiont and nuclear genomes
323 respectively and the ATP biosynthesis cost for the complete biosynthesis of an A:T base pair and a
324 G:C base pair are 40.55 ATP and 40.14 ATP respectively (Chen, et al. 2016). Thus

$$325 \quad \Delta E_{EGT} = A_{DNA} (C_{end} - C_{nuc}) - L_{prot} C_{import} N_p \text{ [9]}$$

326 Where positive values of ΔE_{EGT} correspond to genes for which it is more energetically favourable to
327 be encoded in the nuclear genome, and negative values correspond to genes for which it is more
328 energetically favourable to be encoded in the endosymbiont genome.

329 ***Simulating endosymbiotic gene transfer of mitochondrial and chloroplast genes***

330 The complete genomes with measured protein abundances for an alphaproteobacterium
331 (*Bartonella henselae*) and a cyanobacterium (*Microcystis aeruginosa*) were selected to serve as
332 models for an ancestral mitochondrion and cyanobacterium, respectively. To account for
333 uncertainty in the size and complexity of the ancestral pre-mitochondrial and pre-chloroplast host
334 cells, a range of potential ancestral cells was considered to be engulfed by a range of different host

335 cells with protein contents representative of the diversity of extant eukaryotic cells (Milo 2013).
336 Specifically, the size of the host cell ranged from a small unicellular yeast-like cell (10^7 proteins), to
337 a medium sized unicellular algal-like cell (10^8 proteins) to a typical metazoan/plant cell (10^9
338 proteins). Each of these host cell types was then considered to allocate a realistic range of total
339 cellular protein to mitochondria/chloroplasts typical of eukaryotic cells (i.e. ~2% for yeast (Uchida,
340 et al. 2011), ~20% for metazoan cells (David 1977) and ~50% of the non-vacuolar volume of plant
341 cells (Winter, et al. 1994)). It is not important whether the organellar fraction of the cell is
342 composed of a single large organelle or multiple smaller organelles as all costs, abundances, and
343 copy numbers are evaluated at a per-cell level. For each simulated cell, ΔE_{EGT} was evaluated for
344 each gene in the endosymbiont genome using real protein abundance data (Wang, et al. 2015) for
345 a realistic range of endosymbiont genome copy numbers using equation 9. In all cases the host
346 cell was assumed to be diploid. The simulations were repeated for three different per-residue
347 protein import costs (0.05 ATP, 2 ATP, and 5 ATP per residue respectively). The number of genes
348 where ΔE_{EGT} was positive was recorded as these genes comprise the cohort that are energetically
349 favourable to be encoded in the nuclear genome. All calculated values for ΔE_{EGT} for both the model
350 organisms are provided in Supplemental Table S2.

351 ***Estimating the strength of selection acting on endosymbiotic gene transfer***

352 To model the proportion of energy that would be saved by an individual endosymbiotic gene
353 transfer event a number of assumptions were made. It was assumed that the ancestral host cell
354 had a cell size that is within the range of extant eukaryotes (i.e. between 1×10^7 proteins per cell
355 and 1×10^9 proteins per cell). It was assumed that the endosymbiont occupied a fraction of the
356 total cell proteome that is within the range exhibited by most eukaryotes today (2% to 50% of total
357 cellular protein is located within the endosymbiont under consideration). It was assumed that
358 endosymbiont genome copy number ranged between 1 copy per cell (as it most likely started out
359 with a single copy) and 10,000 copies per cell.

360 We assumed an ancestral host cell with a 24-hour doubling time such that all genomes and
361 proteins are produced in the required abundance every 24-hour period. All cells, irrespective of
362 whether they are bacterial or eukaryotic, consume ATP (C_{ATP}) in proportion to their cell volume (V)
363 at the rate of

364
$$C_{ATP} = 0.39V^{0.88} [10]$$

365 where C_M is in units of 10_9 molecules of ATP cell⁻¹ hour⁻¹, and V is in units of μm^3 (Lynch and
366 Marinov 2015). Thus, the total energy (E_R) needed to replicate a cell was considered to be

367
$$E_R = 24 C_{ATP} [11]$$

368 The proportional energetic advantage or disadvantage ($E_{A/D}$) to the host cell from the
369 endosymbiotic gene transfer of a given gene is evaluated as the free energy of endosymbiotic
370 gene transfer divided by the total amount of energy consumed by the cell during its 24-hour life
371 cycle.

372
$$E_{A/D} = \frac{\Delta E_{EGT}}{E_R} [12]$$

373 Given that $E_{A/D}$ describes the proportional energetic advantage or disadvantage a cell has from a
374 given endosymbiotic gene transfer event $E_{A/D}$ can be used directly as selection coefficient (s) to
375 evaluate the strength of selection acting on the endosymbiotic gene transfer of a given gene. Such
376 that

377
$$s = E_{A/D} [13]$$

378 As ΔE_{EGT} can be positive or negative as described above, s is therefore also positive or negative
379 depending on endosymbiont genome copy number, endosymbiont fraction, host cell protein
380 content, the abundance of the protein that must be imported and the ATP cost of protein import.
381 When s is less than zero the absolute value of s is taken to be the selection coefficient for retention
382 of a gene in the endosymbiont genome (S_R), when s is greater than 0 the value of s is taken to be
383 the selection coefficient for endosymbiotic gene transfer to the nucleus (S_{EGT}). All calculated values
384 for s for both the model alphaproteobacterium (*Bartonella henselae*) and cyanobacterium
385 (*Microcystis aeruginosa*) are provided in Supplemental Table S1.

386 **Estimating time to fixation**

387 Fixation times for endosymbiotic gene transfer events for a range of observed selection coefficients
388 from 1×10^{-5} to 1×10^{-2} were estimated using a Wright–Fisher model with selection and drift
389 (Fisher 1930; Wright 1931) implemented in a simple evolutionary dynamics simulation (Niklaus and
390 Kelly 2018). The effective population size for these simulations was set as 1×10^7 , as is

391 representative of unicellular eukaryotes (Lynch and Conery 2003) and multicellularity in eukaryotes
392 is not thought to have evolved until after the endosymbiosis of either the mitochondrion or the
393 chloroplast.

394 ***The cost of transferring the *rbcL* gene encoding RuBisCO large subunit from the***
395 ***chloroplast to the nuclear genome in *Arabidopsis thaliana****

396 The total number of proteins contained in an *Arabidopsis thaliana* leaf cell is 2.5×10^{10} proteins
397 (Heinemann, et al. 2020). The fraction of cellular protein that is invested in RuBisCO large subunit
398 (F_{rbcL}) is 0.165 (Li, et al. 2017) and thus the number of RuBisCO large subunit proteins per cell (N_p)
399 is estimated to be 4.13×10^9 . The cost of import (P_{import}) of a protein to the chloroplast is 2 ATP per
400 amino acid residue (Shi and Theg 2013). The length of the polypeptide (L_{prot}) comprising the
401 RuBisCO large subunit is 480 amino acids (1440 nucleotides). The ATP biosynthesis cost of a
402 single copy of the *rbcL* gene in double stranded DNA (A_{DNA}) is 58132 ATP. The copy number of the
403 chloroplast genome in a typical *Arabidopsis thaliana* leaf cell (C_{end}) is 1500 copies (Zoschke, et al.
404 2007). *Arabidopsis thaliana* is diploid and thus the copy number of the nuclear genome (C_{nuc}) is 2.
405 Thus, using equation 8 above the ATP that would be saved by transferring the DNA encoding the
406 *rbcL* gene from the chloroplast genome to the nuclear genome is evaluated as

407
$$58132 (1500 - 2) = 8.7 \times 10^7 \text{ ATP [14]}$$

408 Using equation 5 above the ATP that would be required to import the RuBisCO large subunit into
409 the chloroplast is thus evaluated as

410
$$480 \times 2 \times 4.13 \times 10^9 = 3.96 \times 10^{12} \text{ ATP [15]}$$

411 And thus

412
$$\Delta E_{EGT} = 8.7 \times 10^7 - 3.96 \times 10^{12} = -3.96 \times 10^{12} \text{ ATP [16]}$$

413 Given that an *Arabidopsis thaliana* leaf mesophyll cell has a volume of $\sim 49 \mu\text{m}^3$ (Ramonell, et al.
414 2001) the energy consumption rate cell was calculated using equation 10 be consumes 5.2×10^{12}
415 ATP hour^{-1} . Assuming an experimentally determined *in vivo* degradation rate for RuBisCO large
416 subunit of $K_D = 0.052$ (Li, et al. 2017), the recurring daily cost of importing new RuBisCO large
417 subunits into the chloroplast is evaluated as

418
$$L_{prot} C_{import} N_p K_D = 2.1 \times 10^{11} \text{ ATP [17]}$$

419 Thus, if the *rbcL* gene was transferred from the chloroplast genome to the nuclear genome in
420 *Arabidopsis thaliana*, the daily cost of importing the RuBisCO large subunit back into the
421 chloroplast would consume ~0.17% of the total operational energy budget of the cell.

422 **Acknowledgements**

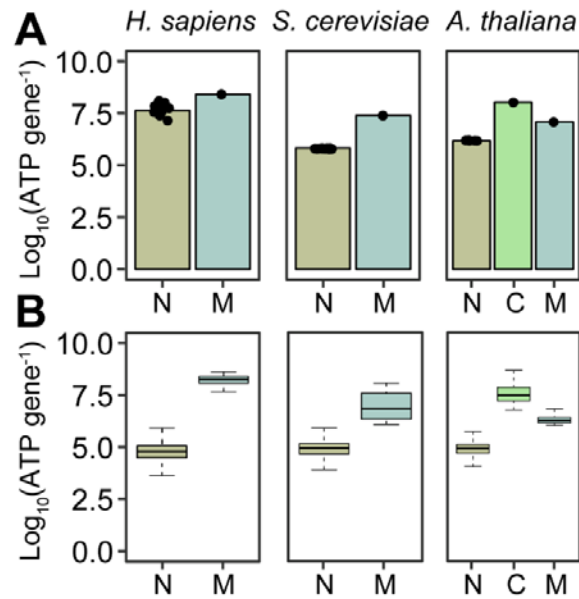
423 This work was funded by the Royal Society and the European Union's Horizon 2020 research and
424 innovation program under grant agreement number 637765. The author would like to thank
425 Thomas A. Richards and John M. Archibald for their comments on the manuscript.

426 **Author Contributions**

427 SK conceived study, conducted the analysis, and wrote the manuscript.

428 **Figures**

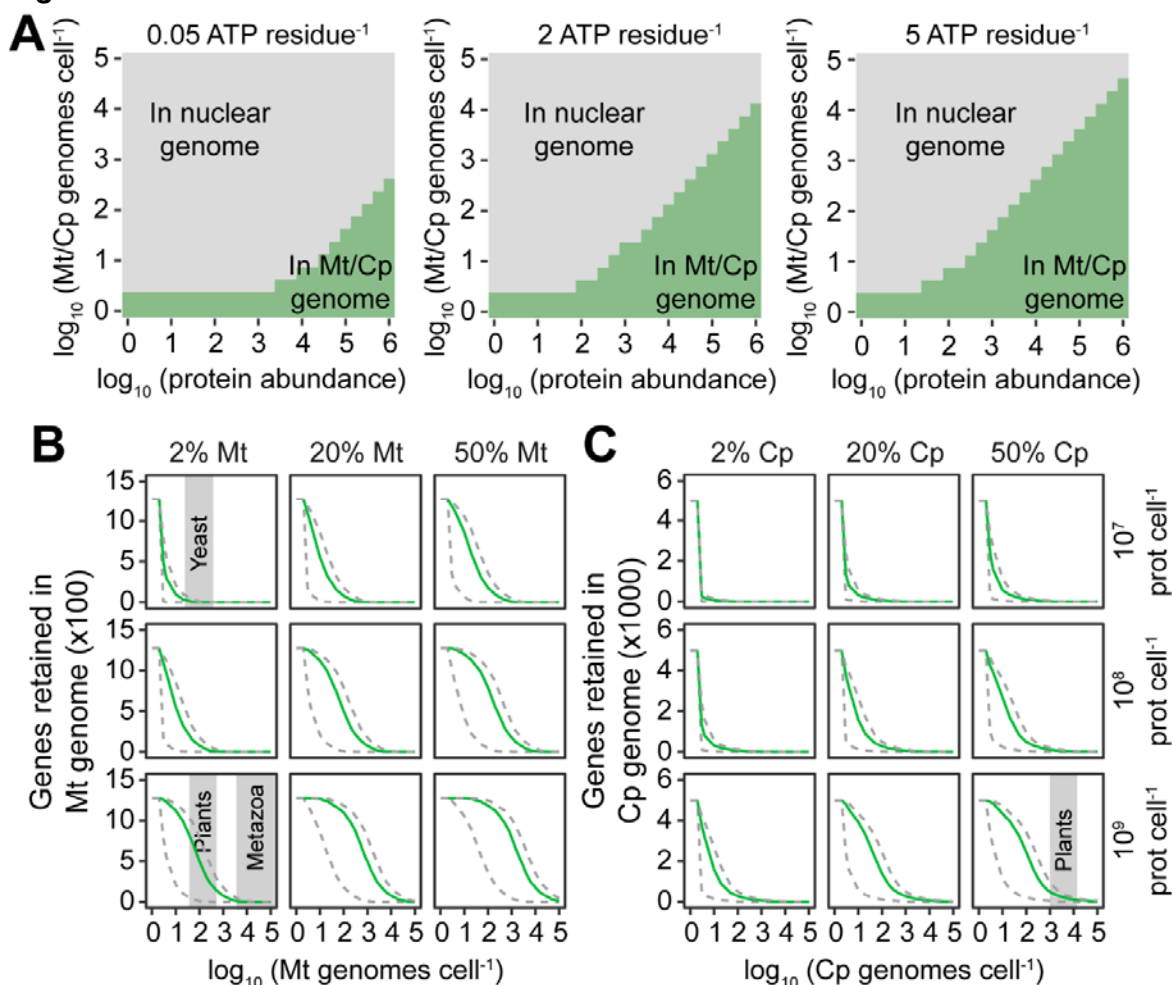
429 **Figure 1**



430

431 **Figure 1.** The per-cell biosynthetic cost of nuclear and organellar genes in three representative
432 eukaryotes. **A)** The ATP biosynthesis costs of nuclear (N), chloroplast (C), and mitochondrial (M)
433 genes calculated as the cost of the chromosome divided by the number of genes contained within
434 that chromosome. Nuclear chromosomes include the cost of nucleosomes, organellar
435 chromosomes only included the cost of the DNA. In the case of the nuclear genes the height of bar
436 depicts the mean cost of all nuclear chromosomes with individual points showing all chromosomes
437 overlaid on top the bar plots. **B)** The ATP biosynthesis cost of just the coding sequences of the
438 genes. In both A and B, the costs were computed assuming a diploid nuclear genome, a per-cell
439 mitochondrial genome copy number of 5000, 200 and 100 for the in *H. sapiens*, *S. cerevisiae* and
440 *A. thaliana*, respectively, and a per cell chloroplast genome copy number of 1500 in *A. thaliana*.

441 **Figure 2**



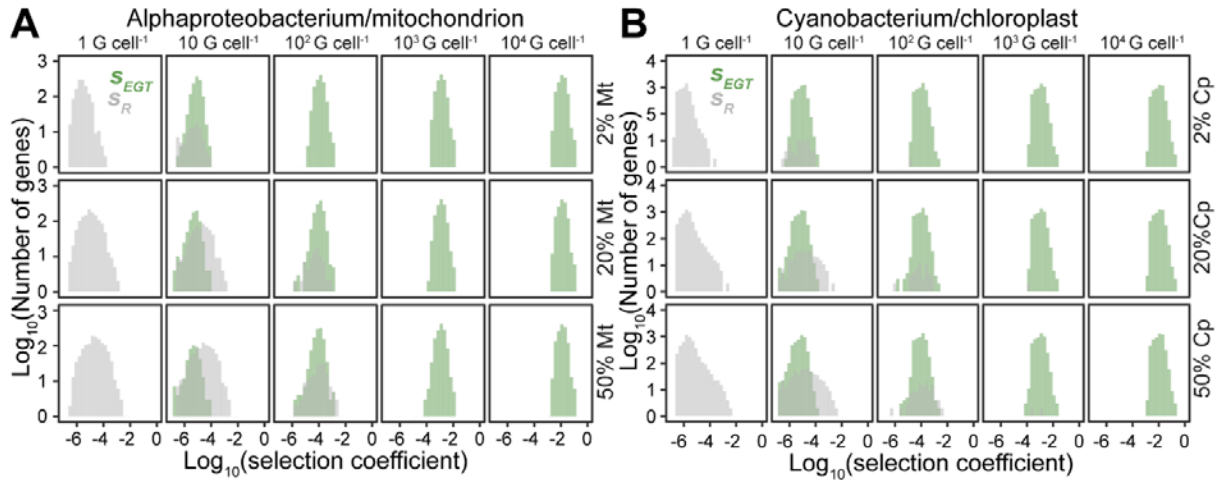
442

443 **Figure 2.** The minimum cost location to the cell of organellar genes encoding an organellar
 444 localised protein. **A)** The minimum cost location of an organellar gene for a range of per-protein
 445 import costs, organellar genome copy numbers, and encoded protein abundance. The grey shaded
 446 fractions of the plots indicate the regions of parameter space where it is more energetically
 447 favourable to the cell to encode an organellar gene in the nuclear genome and import the requisite
 448 amount of protein. The green shaded fractions of the plots indicate the regions of parameter space
 449 where it is more energetically favourable to the cell to encode the gene in the organellar genome.
 450 **B)** The number of genes in the alphaproteobacterial (mitochondrial) genome for which it is more
 451 energetically favourable to the cell for the gene to be retained in the organellar genome. Green
 452 lines assume a per-residue protein import cost of 2 ATP per amino acid. Grey dashed lines
 453 indicate lower and upper cost bounds of 0.05 ATP and 5 ATP per residue respectively. **C)** As in B
 454 but for the cyanobacterial (chloroplast) genome. Grey shaded areas on plots are provided for

- 455 illustrative purposes to indicate the regions of parameter space occupied by yeast, metazoan and
456 plant cells. Cp: chloroplast. Mt: mitochondrion.

457 **Figure 3**

458



459 **Figure 3.** Selection coefficients for retention (S_R , grey) or endosymbiotic gene transfer (S_{EGT} ,
 460 green) of all genes encoded in the example alphaproteobacterial and cyanobacterial genomes.
 461 Coefficients were computed accounting for protein abundance, host cell organellar fraction,
 462 organellar genome copy number per cell, and host cell energy consumption. Plots shown are for a
 463 simulated host cell comprising 1×10^7 proteins and a protein import cost of 2 ATP per residue,
 464 plots for other host cell protein contents and protein import costs are provided in Supplemental
 465 Figures S1-S6. **A)** Selection coefficients of all genes encoded in the alphaproteobacterium
 466 genome. **B)** Selection coefficients for all genes encoded in the cyanobacterial genome. S_R and
 467 S_{EGT} have opposite signs (see methods). To simplify the display and enable direct comparison, the
 468 absolute value of the selection coefficients of each gene are plotted and green shading is used to
 469 indicate genes in the S_{EGT} fraction and grey shading indicates genes in the S_R fraction of the
 470 genome. Mt, mitochondrion. Cp, chloroplast. G, genomes.

471 **References**

- 472 Allen JF, Raven JA. 1996. Free-radical-induced mutation vs redox regulation:
473 costs and benefits of genes in organelles. *J Mol Evol* 42:482-492.
- 474 Almagro Armenteros JJ, Salvatore M, Emanuelsson O, Winther O, von Heijne
475 G, Elofsson A, Nielsen H. 2019. Detecting sequence signals in targeting
476 peptides using deep learning. *Life Sci Alliance* 2.
- 477 Archibald John M. 2015a. Endosymbiosis and Eukaryotic Cell Evolution.
478 *Current Biology* 25:R911-R921.
- 479 Archibald JM. 2015b. Genomic perspectives on the birth and spread of
480 plastids. *Proceedings of the National Academy of Sciences* 112:10147-
481 10153.
- 482 Backes S, Herrmann JM. 2017. Protein Translocation into the Intermembrane
483 Space and Matrix of Mitochondria: Mechanisms and Driving Forces. *Frontiers*
484 *in Molecular Biosciences* 4.
- 485 Bar-On YM, Phillips R, Milo R. 2018. The biomass distribution on Earth.
486 *Proceedings of the National Academy of Sciences* 115:6506-6511.
- 487 Bendich AJ. 1987. Why do chloroplasts and mitochondria contain so many
488 copies of their genome? *Bioessays* 6:279-282.
- 489 Bersaglieri T, Sabeti PC, Patterson N, Vanderploeg T, Schaffner SF, Drake
490 JA, Rhodes M, Reich DE, Hirschhorn JN. 2004. Genetic signatures of strong
491 recent positive selection at the lactase gene. *Am J Hum Genet* 74:1111-1120.
- 492 Booth A, Doolittle WF. 2015a. Eukaryogenesis, how special really?
493 *Proceedings of the National Academy of Sciences* 112:10278-10285.
- 494 Booth A, Doolittle WF. 2015b. Reply to Lane and Martin: Being and becoming
495 eukaryotes. *Proceedings of the National Academy of Sciences* 112:E4824-
496 E4824.
- 497 Brown JR. 2003. Ancient horizontal gene transfer. *Nature Reviews Genetics*
498 4:121-132.
- 499 Calvo SE, Mootha VK. 2010. The mitochondrial proteome and human
500 disease. *Annu Rev Genomics Hum Genet* 11:25-44.
- 501 Chen W-H, Lu G, Bork P, Hu S, Lercher MJ. 2016. Energy efficiency trade-
502 offs drive nucleotide usage in transcribed regions. *Nature Communications*
503 7:11334.
- 504 Cherry JM, Hong EL, Amundsen C, Balakrishnan R, Binkley G, Chan ET,
505 Christie KR, Costanzo MC, Dwight SS, Engel SR, et al. 2012.
506 *Saccharomyces Genome Database: the genomics resource of budding yeast.*
507 *Nucleic Acids Res* 40:D700-705.
- 508 Cole LW. 2016. The Evolution of Per-cell Organelle Number. *Front Cell Dev*
509 *Biol* 4:85.
- 510 Dagan T, Roettger M, Stucken K, Landan G, Koch R, Major P, Gould SB,
511 Goremykin VV, Rippka R, Tandeau de Marsac N, et al. 2013. Genomes of
512 Stigonematalean cyanobacteria (subsection V) and the evolution of oxygenic
513 photosynthesis from prokaryotes to plastids. *Genome Biol Evol* 5:31-44.
- 514 Daley DO, Whelan J. 2005. Why genes persist in organelle genomes.
515 *Genome Biology* 6:110.

- 516 David H. 1977. Quantitative Ultrastructural Data of Animal and Human Cells:
517 Gustav Fischer.
- 518 Deusch O, Landan G, Roettger M, Gruenheit N, Kowallik KV, Allen JF, Martin
519 W, Dagan T. 2008. Genes of Cyanobacterial Origin in Plant Nuclear
520 Genomes Point to a Heterocyst-Forming Plastid Ancestor. *Molecular Biology
521 and Evolution* 25:748-761.
- 522 Doolittle WF. 1998. You are what you eat: a gene transfer ratchet could
523 account for bacterial genes in eukaryotic nuclear genomes. *Trends Genet*
524 14:307-311.
- 525 Drouin G, Daoud H, Xia J. 2008. Relative rates of synonymous substitutions
526 in the mitochondrial, chloroplast and nuclear genomes of seed plants. *Mol
527 Phylogenet Evol* 49:827-831.
- 528 Ferro M, Brugière S, Salvi D, Seigneurin-Berny D, Court M, Moyet L, Ramus
529 C, Miras S, Mellal M, Le Gall S, et al. 2010. AT_CHLORO, a comprehensive
530 chloroplast proteome database with subplastidial localization and curated
531 information on envelope proteins. *Mol Cell Proteomics* 9:1063-1084.
- 532 Fisher RAS. 1930. The genetical theory of natural selection. Oxford:
533 Clarendon Press.
- 534 Goodstein DM, Shu S, Howson R, Neupane R, Hayes RD, Fazo J, Mitros T,
535 Dirks W, Hellsten U, Putnam N, et al. 2012. Phytozome: a comparative
536 platform for green plant genomics. *Nucleic Acids Res* 40:D1178-1186.
- 537 Gray MW, Burger G, Lang BF. 1999. Mitochondrial evolution. *Science*
538 283:1476-1481.
- 539 Green BR. 2011. Chloroplast genomes of photosynthetic eukaryotes. *Plant J*
540 66:34-44.
- 541 Grisdale CJ, Smith DR, Archibald JM. 2019. Relative Mutation Rates in
542 Nucleomorph-Bearing Algae. *Genome Biology and Evolution* 11:1045-1053.
- 543 Hartl DL, Moriyama EN, Sawyer SA. 1994. Selection intensity for codon bias.
544 *Genetics* 138:227-234.
- 545 Heinemann B, Künzler P, Braun H-P, Hildebrandt TM. 2020. Estimating the
546 number of protein molecules in a plant cell: a quantitative perspective on
547 proteostasis and amino acid homeostasis during progressive drought stress.
548 *bioRxiv:2020.2003.2017.995613*.
- 549 Herrmann R. 1997. Eukaryotism, towards a new interpretation. In.
550 *Eukaryotism and symbiosis: Springer*. p. 73-118.
- 551 Husnik F, Nikoh N, Koga R, Ross L, Duncan RP, Fujie M, Tanaka M, Satoh
552 N, Bachtrog D, Wilson AC, et al. 2013. Horizontal gene transfer from diverse
553 bacteria to an insect genome enables a tripartite nested mealybug symbiosis.
554 *Cell* 153:1567-1578.
- 555 Jarvis P. 2008. Targeting of nucleus-encoded proteins to chloroplasts in
556 plants. *New Phytologist* 179:257-285.
- 557 Khakhlova O, Bock R. 2006. Elimination of deleterious mutations in plastid
558 genomes by gene conversion. *The Plant Journal* 46:85-94.
- 559 Lane N. 2014. Bioenergetic constraints on the evolution of complex life. *Cold
560 Spring Harb Perspect Biol* 6:a015982.

- 561 Lane N, Martin W. 2010. The energetics of genome complexity. *Nature*
562 467:929-934.
- 563 Lane N, Martin WF. 2015. Eukaryotes really are special, and mitochondria
564 are why. *Proceedings of the National Academy of Sciences* 112:E4823-
565 E4823.
- 566 Li L, Nelson CJ, Trösch J, Castleden I, Huang S, Millar AH. 2017. Protein
567 Degradation Rate in *Arabidopsis thaliana* Leaf Growth and
568 Development. *The Plant Cell* 29:207-228.
- 569 Lynch M. 1997. Mutation accumulation in nuclear, organelle, and prokaryotic
570 transfer RNA genes. *Mol Biol Evol* 14:914-925.
- 571 Lynch M. 1996. Mutation accumulation in transfer RNAs: molecular evidence
572 for Muller's ratchet in mitochondrial genomes. *Mol Biol Evol* 13:209-220.
- 573 Lynch M, Conery JS. 2003. The origins of genome complexity. *Science*
574 302:1401-1404.
- 575 Lynch M, Koskella B, Schaack S. 2006. Mutation Pressure and the Evolution
576 of Organelle Genomic Architecture. *Science* 311:1727-1730.
- 577 Lynch M, Lynch PSTSM, Walsh B. 2007. *The Origins of Genome*
578 *Architecture*: Oxford University Press, Incorporated.
- 579 Lynch M, Marinov GK. 2015. The bioenergetic costs of a gene. *Proceedings*
580 *of the National Academy of Sciences* 112:15690-15695.
- 581 Lynch M, Marinov GK. 2017. Membranes, energetics, and evolution across
582 the prokaryote-eukaryote divide. *eLife* 6:e20437.
- 583 Lynch M, Marinov GK. 2018. Response to Martin and colleagues:
584 mitochondria do not boost the bioenergetic capacity of eukaryotic cells.
585 *Biology Direct* 13:26.
- 586 Martin W, Herrmann RG. 1998. Gene Transfer from Organelles to the
587 Nucleus: How Much, What Happens, and Why? *Plant Physiology* 118:9-17.
- 588 Martin W, Kowallik K. 1999. Annotated English translation of
589 Mereschkowsky's 1905 paper 'Über Natur und Ursprung der Chromatophoren
590 im Pflanzenreiche'. *European Journal of Phycology* 34:287-295.
- 591 Martin W, Müller M. 1998. The hydrogen hypothesis for the first eukaryote.
592 *Nature* 392:37-41.
- 593 Martin W, Rujan T, Richly E, Hansen A, Cornelsen S, Lins T, Leister D,
594 Stoebe B, Hasegawa M, Penny D. 2002. Evolutionary analysis of
595 *Arabidopsis*, cyanobacterial, and chloroplast genomes reveals
596 plastid phylogeny and thousands of cyanobacterial genes in the nucleus.
597 *Proceedings of the National Academy of Sciences* 99:12246-12251.
- 598 Martin WF, Garg S, Zimorski V. 2015. Endosymbiotic theories for eukaryote
599 origin. *Philos Trans R Soc Lond B Biol Sci* 370:20140330.
- 600 McCutcheon JP, Moran NA. 2012. Extreme genome reduction in symbiotic
601 bacteria. *Nature Reviews Microbiology* 10:13-26.
- 602 Milo R. 2013. What is the total number of protein molecules per cell volume?
603 A call to rethink some published values. *Bioessays* 35:1050-1055.
- 604 Miyakawa I. 2017. Organization and dynamics of yeast mitochondrial
605 nucleoids. *Proc Jpn Acad Ser B Phys Biol Sci* 93:339-359.

- 606 Mokranjac D, Neupert W. 2008. Energetics of protein translocation into
607 mitochondria. *Biochimica et Biophysica Acta (BBA) - Bioenergetics* 1777:758-
608 762.
- 609 Muller HJ. 1964. The relation of recombination to mutational advance.
610 *Mutation Research/Fundamental and Molecular Mechanisms of Mutagenesis*
611 1:2-9.
- 612 Nakayama T, Ishida K. 2009. Another acquisition of a primary photosynthetic
613 organelle is underway in *Paulinella chromatophora*. *Curr Biol* 19:R284-285.
- 614 Neiman M, Taylor DR. 2009. The causes of mutation accumulation in
615 mitochondrial genomes. *Proceedings of the Royal Society B: Biological*
616 *Sciences* 276:1201-1209.
- 617 Niklaus M, Kelly S. 2018. The molecular evolution of C4 photosynthesis:
618 opportunities for understanding and improving the world's most productive
619 plants. *Journal of Experimental Botany* 70:795-804.
- 620 Nowack ECM, Vogel H, Groth M, Grossman AR, Melkonian M, Glöckner G.
621 2010. Endosymbiotic Gene Transfer and Transcriptional Regulation of
622 Transferred Genes in *Paulinella chromatophora*. *Molecular Biology and*
623 *Evolution* 28:407-422.
- 624 Nowack ECM, Weber APM. 2018. Genomics-Informed Insights into
625 Endosymbiotic Organelle Evolution in Photosynthetic Eukaryotes. *Annual*
626 *Review of Plant Biology* 69:51-84.
- 627 Ramonell KM, Kuang A, Porterfield DM, Crispi ML, Xiao Y, McClure G,
628 Musgrave ME. 2001. Influence of atmospheric oxygen on leaf structure and
629 starch deposition in *Arabidopsis thaliana*. *Plant Cell Environ* 24:419-428.
- 630 Reyes-Prieto A. 2015. The basic genetic toolkit to move in with your
631 photosynthetic partner. *Frontiers in Ecology and Evolution* 3.
- 632 Reyes-Prieto A, Hackett JD, Soares MB, Bonaldo MF, Bhattacharya D. 2006.
633 Cyanobacterial Contribution to Algal Nuclear Genomes Is Primarily Limited to
634 Plastid Functions. *Current Biology* 16:2320-2325.
- 635 Reyes-Prieto A, Yoon HS, Moustafa A, Yang EC, Andersen RA, Boo SM,
636 Nakayama T, Ishida K, Bhattacharya D. 2010. Differential gene retention in
637 plastids of common recent origin. *Mol Biol Evol* 27:1530-1537.
- 638 Roger AJ, Muñoz-Gómez SA, Kamikawa R. 2017. The Origin and
639 Diversification of Mitochondria. *Curr Biol* 27:R1177-r1192.
- 640 Rogozin IB, Carmel L, Csuros M, Koonin EV. 2012. Origin and evolution of
641 spliceosomal introns. *Biol Direct* 7:11.
- 642 Schatz G, Dobberstein B. 1996. Common Principles of Protein Translocation
643 Across Membranes. *Science* 271:1519-1526.
- 644 Shi L-X, Theg SM. 2013. Energetic cost of protein import across the envelope
645 membranes of chloroplasts. *Proceedings of the National Academy of*
646 *Sciences* 110:930-935.
- 647 Shokolenko I, Venediktova N, Bochkareva A, Wilson GL, Alexeyev MF. 2009.
648 Oxidative stress induces degradation of mitochondrial DNA. *Nucleic Acids*
649 *Res* 37:2539-2548.

- 650 Singer A, Poschmann G, Mühlich C, Valadez-Cano C, Hänsch S, Hüren V,
651 Rensing SA, Stühler K, Nowack ECM. 2017. Massive Protein Import into the
652 Early-Evolutionary-Stage Photosynthetic Organelle of the Amoeba *Paulinella*
653 chromatophora. *Curr Biol* 27:2763-2773.e2765.
- 654 Smith DR. 2015. Mutation rates in plastid genomes: they are lower than you
655 might think. *Genome Biol Evol* 7:1227-1234.
- 656 Smith DR. 2016. The mutational hazard hypothesis of organelle genome
657 evolution: 10 years on. *Molecular Ecology* 25:3769-3775.
- 658 Smith DR, Keeling PJ. 2015. Mitochondrial and plastid genome architecture:
659 Reoccurring themes, but significant differences at the extremes. *Proceedings*
660 *of the National Academy of Sciences* 112:10177-10184.
- 661 Soll J, Schleiff E. 2004. Protein import into chloroplasts. *Nat Rev Mol Cell Biol*
662 5:198-208.
- 663 Speijer D, Hammond M, Lukeš J. 2020. Comparing Early Eukaryotic
664 Integration of Mitochondria and Chloroplasts in the Light of Internal ROS
665 Challenges: Timing is of the Essence. *mBio* 11:e00955-00920.
- 666 Thiergart T, Landan G, Schenk M, Dagan T, Martin WF. 2012. An
667 evolutionary network of genes present in the eukaryote common ancestor
668 polls genomes on eukaryotic and mitochondrial origin. *Genome Biol Evol*
669 4:466-485.
- 670 Timmis JN, Ayliffe MA, Huang CY, Martin W. 2004. Endosymbiotic gene
671 transfer: organelle genomes forge eukaryotic chromosomes. *Nat Rev Genet*
672 5:123-135.
- 673 Uchida M, Sun Y, McDermott G, Knoechel C, Le Gros MA, Parkinson D,
674 Drubin DG, Larabell CA. 2011. Quantitative analysis of yeast internal
675 architecture using soft X-ray tomography. *Yeast* 28:227-236.
- 676 Wang M, Herrmann CJ, Simonovic M, Szklarczyk D, von Mering C. 2015.
677 Version 4.0 of PaxDb: Protein abundance data, integrated across model
678 organisms, tissues, and cell-lines. *Proteomics* 15:3163-3168.
- 679 Wiedemann N, Pfanner N. 2017. Mitochondrial Machineries for Protein Import
680 and Assembly. *Annual Review of Biochemistry* 86:685-714.
- 681 Winter H, Robinson DG, Heldt HW. 1994. Subcellular volumes and metabolite
682 concentrations in spinach leaves. *Planta* 193:530-535.
- 683 Wolfe KH, Li WH, Sharp PM. 1987. Rates of nucleotide substitution vary
684 greatly among plant mitochondrial, chloroplast, and nuclear DNAs.
685 *Proceedings of the National Academy of Sciences* 84:9054-9058.
- 686 Wright S. 1931. Evolution in Mendelian Populations. *Genetics* 16:97-159.
- 687 Yang D, Oyaizu Y, Oyaizu H, Olsen GJ, Woese CR. 1985. Mitochondrial
688 origins. *Proceedings of the National Academy of Sciences* 82:4443-4447.
- 689 Yates AD, Achuthan P, Akanni W, Allen J, Allen J, Alvarez-Jarreta J, Amode
690 MR, Armean IM, Azov AG, Bennett R, et al. 2020. Ensembl 2020. *Nucleic*
691 *Acids Res* 48:D682-d688.
- 692 Zoschke R, Liere K, Börner T. 2007. From seedling to mature plant:
693 *Arabidopsis* plastidial genome copy number, RNA accumulation and

694 transcription are differentially regulated during leaf development. The Plant
695 Journal 50:710-722.
696