

Journal Pre-proof

Conservation and over-representation of G-quadruplex sequences in regulatory regions of mitochondrial DNA across distinct taxonomic sub-groups

Natália Bohálová, Michaela Dobrovolná, Václav Brázda, Stefan Bidula



PII: S0300-9084(21)00294-7

DOI: <https://doi.org/10.1016/j.biochi.2021.12.006>

Reference: BIOCHI 6227

To appear in: *Biochimie*

Received Date: 8 October 2021

Revised Date: 22 November 2021

Accepted Date: 14 December 2021

Please cite this article as: Natá. Bohálová, M. Dobrovolná, Vá. Brázda, S. Bidula, Conservation and over-representation of G-quadruplex sequences in regulatory regions of mitochondrial DNA across distinct taxonomic sub-groups, *Biochimie* (2022), doi: <https://doi.org/10.1016/j.biochi.2021.12.006>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2021 Published by Elsevier B.V.

Author contributions

Natália Bohálová: Software, Validation, Formal analysis, Investigation, Writing – Review and Editing, Visualization

Michaela Dobrovolná: Formal analysis, Investigation, Resources, Writing – Review and Editing

Václav Brázda: Conceptualization, Resources, Writing - Review and Editing, Supervision, Project Administration, Funding acquisition

Stefan Bidula: Conceptualization, Formal analysis, Investigation, Writing - Original Draft, Visualization, Supervision, Project Administration

Abstract

G-quadruplexes have important regulatory roles in the nuclear genome but their distribution and potential roles in mitochondrial DNA (mtDNA) are poorly understood. We analysed 11883 mtDNA sequences from 18 taxonomic sub-groups and identified their frequency and location within mtDNA. Large differences in both the frequency and number of putative quadruplex-forming sequences (PQS) were observed amongst all the organisms and PQS frequency was negatively correlated with an increase in evolutionary age. PQS were over-represented in the 3'UTRs, D-loops, replication origins, and stem loops, indicating regulatory roles for quadruplexes in mtDNA. Variations of the G-quadruplex-forming sequence in the conserved sequence block II (CSBII) region of the human D-loop were conserved amongst other mammals, amphibians, birds, reptiles, and fishes. This D-loop PQS was conserved in the duplicated control regions of some birds and reptiles, indicating its importance to mitochondrial function. The guanine tracts in these PQS also displayed significant length heterogeneity and the length of these guanine tracts were generally longest in bird mtDNA. This information provides further insights into how G4s may contribute to the regulation and function of mtDNA and acts as a database of information for future studies investigating mitochondrial G4s in organisms other than humans.

Key words: mitochondria, genome, G-quadruplex, evolution, D-loop

1 **Conservation and over-representation of G-quadruplex sequences in regulatory**
2 **regions of mitochondrial DNA across distinct taxonomic sub-groups**

3 Running title – G-quadruplexes in mitochondria

4

5 Natália Bohálová,^{a, b} Michaela Dobrovolná,^{a, c} Václav Brázda,^{a, c} Stefan Bidula^d

6

7 ^a Institute of Biophysics of the Czech Academy of Sciences, Brno, Czech Republic

8 ^b Department of Experimental Biology, Faculty of Science, Masaryk University, Brno, Czech
9 Republic

10 ^c Department of Food Chemistry and Biotechnology, Faculty of Chemistry, Brno University of
11 Technology, Purkyňova 118, 61200 Brno, Czech Republic

12 ^d School of Biological Sciences, University of East Anglia, Norwich, UK

13

14 Corresponding author – Stefan Bidula (s.bidula@uea.ac.uk), University of East Anglia,
15 Norwich Research Park, Norwich, NR4 7TJ

16

17 Declarations of interest: None

18

19 **Abbreviations**

20 G4, G-quadruplex; PQS. Putative quadruplex-forming sequence; mitochondrial DNA, mtDNA;

21 CBSII, conserved sequence block II; RHPS4, 3,11-Difluoro-6,8,13-trimethylquino[4,3,2-

22 *k*]acridinium methylsulfate; NCBI, National Center for Biotechnology and Information

23

24 **Abstract**

25 G-quadruplexes have important regulatory roles in the nuclear genome but their distribution
26 and potential roles in mitochondrial DNA (mtDNA) are poorly understood. We analysed 11883
27 mtDNA sequences from 18 taxonomic sub-groups and identified their frequency and location
28 within mtDNA. Large differences in both the frequency and number of putative quadruplex-
29 forming sequences (PQS) were observed amongst all the organisms and PQS frequency was
30 negatively correlated with an increase in evolutionary age. PQS were over-represented in the
31 3'UTRs, D-loops, replication origins, and stem loops, indicating regulatory roles for
32 quadruplexes in mtDNA. Variations of the G-quadruplex-forming sequence in the conserved
33 sequence block II (CSBII) region of the human D-loop were conserved amongst other
34 mammals, amphibians, birds, reptiles, and fishes. This D-loop PQS was conserved in the
35 duplicated control regions of some birds and reptiles, indicating its importance to mitochondrial
36 function. The guanine tracts in these PQS also displayed significant length heterogeneity and
37 the length of these guanine tracts were generally longest in bird mtDNA. This information
38 provides further insights into how G4s may contribute to the regulation and function of mtDNA
39 and acts as a database of information for future studies investigating mitochondrial G4s in
40 organisms other than humans.

41

42 **Key words:** mitochondria, genome, G-quadruplex, evolution, D-loop

43

44

45

46

47

48

49 1. Introduction

50 Mitochondria are aptly referred to as the 'powerhouse' of the cell, whereby they generate
51 energy for cells in the form of ATP [1]. However, they have now been demonstrated to
52 participate in numerous important and diverse biological processes, including metabolic
53 signalling, bioenergetics, calcium transport, production of reactive oxygen species, and
54 regulation of cell death pathways [2]. Dysfunction of mitochondria, arising from either the
55 acquisition of mutations in mitochondrial DNA (mtDNA) or nuclear DNA, can be catastrophic
56 and can result in various diseases, such as Leigh syndrome, myoclonic epilepsy with ragged
57 red fibres syndrome, and mtDNA depletion syndrome [3]. Thus, mitochondria require strict
58 regulatory processes to ensure normal biological function. There is growing evidence that
59 alternative DNA structures such as cruciforms, left-handed DNA (Z-DNA), R-loops, and
60 quadruplexes play critical regulatory roles in fundamental biological functions, although our
61 understanding of their roles in mtDNA are still in their infancy [4–7].

62 G-quadruplexes (G4s) are four-stranded secondary structures in nucleic acids that form in
63 guanine-rich regions. G4s form when four guanines associate through Hoogsteen hydrogen
64 bonding to form a G-tetrad [8]. Several G-tetrads then stack on top of one another, linked by
65 mixed-sequence nucleotides, to form the G4 structure itself. Given their localisation
66 throughout the nuclear genome in promoters, untranslated regions, and telomeres, G4s and
67 iMs have been identified to participate in critical regulatory processes, such as transcription,
68 translation, and phase separation of RNA to name but a few [9,10]. However, their roles within
69 mtDNA are poorly understood and practically nothing known about their roles in the mtDNA of
70 organisms other than humans.

71

72 Therefore, it was important that we explored mtDNA for the presence of putative quadruplex-
73 forming sequences (PQS) to enhance our understanding of where these structures were
74 located in mtDNA and provide insight into the potential biological roles that G4s may play. To

75 this end, we have analysed the mtDNA of 11883 genomes from diverse species across 18
76 taxonomic sub-groups and highlighted the frequency and location of PQS, whilst also
77 identifying conservation of a PQS in a critical regulatory region across taxonomic sub-groups.

78

79 **2. Methods**

80 **2.1 mtDNA sequences**

81 Complete and most recent mtDNA sequences were downloaded from the organelle genome
82 database of the National Center for Biotechnology and Information (NCBI). The genomes were
83 obtained for 11883 organisms from 18 different taxonomic sub-groups. Groups included
84 amphibians (303 genomes), apicomplexans (47 genomes), ascomycetes (379 genomes),
85 basidiomycetes (125 genomes), birds (977 genomes), fishes (3036 genomes), flatworms (153
86 genomes), green algae (92 genomes), insects (2534 genomes), land plants (259 genomes),
87 mammals (1342 genomes), other (161 genomes), other animals (1760 genomes), other fungi
88 (27 genomes), other plants (13 genomes), other protists (109 genomes), reptiles (382
89 genomes), and roundworms (184 genomes). Duplicated genomes were omitted from the
90 analysis.

91

92 **2.2 Data analysis**

93 Genomes were analysed using G4Hunter (<http://bioinformatics.ibp.cz>) to identify PQS [11].
94 The parameters for analysis were set at a length of 25 nucleotides and a threshold of 1.2, as
95 these settings have previously been shown to identify experimentally validated quadruplex
96 structures [12]. The figures in the main manuscript body were all produced using these
97 analysis parameters. However, we also analysed the genomes with thresholds between 1.2-
98 1.4, 1.4-1.6, 1.6-1.8, 1.8-2.0, and more than 2.0. The raw data for all the genomes analysed
99 can be found in **Supplementary Table S1** and provides information on the genome name,

100 NCBI identifier, length, GC genome content, and frequency and number of PQS. The tables
101 of annotated genomic features were downloaded in tandem with the mitochondrial genomes,
102 and we analysed the frequency of PQS within these annotated features (e.g., gene) and within
103 ± 100 base pairs of the features for each genome (**Supplementary Table S2**). The script used
104 for analysis is publicly available at https://gitlab.com/PatrikKaura/DNA_analyser_IBP.

105

106 **2.3 Statistical analysis**

107 A cluster dendrogram was constructed in R, using the *pvclust* package [13]. The following
108 values were used as input data: Mean PQS/kbp, Min PQS/kbp, Max PQS/kbp. and Cov % (%
109 of genome covered by PQS). The 'ward.D2' clustering method was used with Euclidean
110 distance and 10,000 bootstrap resampling. This cluster dendrogram can be found in
111 **Supplementary Figure S2**. Correlation was determined by two-tailed Pearson's correlation
112 coefficient. Normality of the data was determined via a Shapiro-Wilk test. Non-parametric
113 Kruskal-Wallis tests with Dunn's multiple comparisons were used to determine significance.
114 All figures and analysis were generated using GraphPad Prism (v 9.1.0).

115

116 **3. Results and discussion**

117 **3.1 Large heterogeneity of PQS frequency in mtDNA**

118 G4s and iMs have been shown to have important regulatory roles throughout nuclear
119 genomes. However, an in-depth analysis of PQS in mitochondrial genomes had not been
120 conducted to date. Using G4Hunter, we investigated the presence and frequency of PQS in
121 11883 mitochondrial sequences from highly varied species across 18 taxonomic sub-groups.

122

123 We utilised the default G4Hunter settings to identify PQS with a length of 25 nt using a
124 threshold of 1.2. Data for all organisms can be found in **Supplementary Table S1** and

125 **Supplementary Figure S1.** The PQS frequencies between all groups were found to be
126 significantly different (**Figure 1**). On average, bird mtDNA was found to have the highest PQS
127 frequency (5.62 PQS/kbp), frequency of PQS relative to genome GC content (12.28
128 PQS/GC%), and highest genome GC content (45.63%; **Figure 1, Table 1**). Conversely, the
129 mtDNA of land plants was found to have the greatest average total number of PQS (607.03
130 PQS; **Table 1**). Conversely, apicomplexan mtDNA contained the lowest PQS/kbp (0.16
131 PQS/kbp), PQS/GC% (0.49 PQS/GC%), and total PQS (0.94 PQS), whilst insect mtDNA
132 contained the lowest genome GC content (23.78%; **Figure 1, Table 1**).

133

134 In general, mtDNA PQS frequencies were associated with evolutionary distance. More closely
135 related organisms, such as birds, reptiles, mammals, amphibians, and fishes had the highest
136 PQS frequencies and clustered together with statistical significance, whereas there appeared
137 to be a loss of PQS with an increase in evolutionary distance (**Supplementary Figure S2**).

138

139 Interestingly, the frequency of PQS was found to be the inverse of what has previously been
140 observed for inverted repeats in mtDNA [14]. Inverted repeats of six or more nucleotides are
141 a pre-requisite for the formation of cruciforms; alternative DNA structures which have been
142 demonstrated to regulate key biological processes [4,15]. Birds were found to have the lowest
143 frequency of inverted repeats, whilst apicomplexans, fungi, and insects had the highest
144 frequency [14]. This is almost a mirror-image of what we observed for PQS and could be
145 indicative that inverted repeats are more important for the integrity or function of mitochondria
146 in organisms such as apicomplexans, whereas inverted repeats may be detrimental to more
147 complex organisms such as birds. Indeed, inverted repeat-related inversions tend to
148 accumulate in tissues with high-energy metabolism in mammals, which can accelerate ageing
149 and reduce longevity [16]. Birds possess fewer inverted repeats and have much longer
150 lifespans comparative to similarly sized organisms [17]. Our understanding of the roles of

151 quadruplexes within mtDNA, particularly in organisms other than humans, are still in their
 152 infancy. Although, one could speculate that the roles of G4s in mtDNA could be organism
 153 specific and dependent upon the prevalence of PQS within the genome.

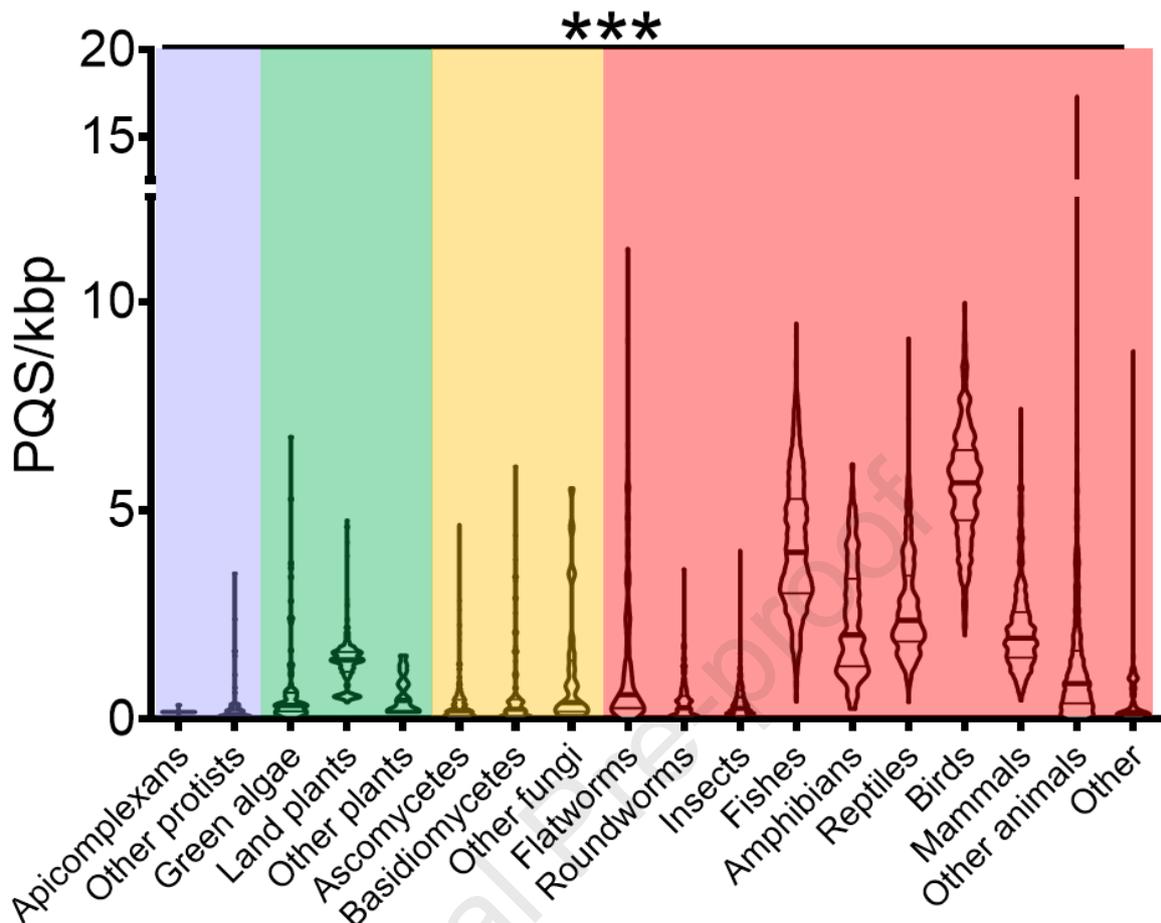
154

155

156 **Table 1.** An overview of the PQS frequencies, PQS number, and genome GC% of
 157 mitochondrial genomes in the study

Group name	Number of genomes	Average PQS/kbp	Lowest frequency (PQS/kbp)	Highest frequency (PQS/kbp)	Average genome GC%	Average PQS/GC%	Average total PQS
Apicomplexans	47	0.16	0.00	0.34	31.62	0.49	0.94
Other protists	109	0.24	0.00	3.49	27.17	0.74	22.65
Green algae	92	0.70	0.03	6.76	37.74	1.58	50.68
Land plants	259	1.38	0.40	4.76	44.10	3.08	607.03
Other plants	13	0.61	0.17	1.52	37.84	1.51	70.15
Ascomycetes	379	0.43	0.00	4.66	26.42	1.62	30.24
Basidiomycetes	125	0.59	0.00	6.05	28.09	1.87	55.35
Other fungi	27	1.14	0.07	5.53	33.97	2.78	87.22
Flatworms	153	1.37	0.00	11.27	31.73	3.76	20.05
Roundworms	184	0.41	0.00	3.59	25.39	1.46	5.89
Insects	2534	0.36	0.00	4.04	23.78	1.41	5.71
Fishes	3036	4.17	0.42	9.49	44.16	9.31	69.61
Amphibians	303	2.38	0.24	6.11	38.87	5.91	42.46
Reptiles	382	2.68	0.40	9.12	40.56	6.49	45.98
Birds	977	5.62	2.00	9.98	45.60	12.28	96.35
Mammals	1342	2.15	0.44	7.44	39.18	5.39	35.71
Other animals	1760	1.32	0.00	17.34	32.70	3.67	22.44
Other	161	0.44	0.00	8.82	31.35	1.24	14.63

158



159

160 **Figure 1. There was large heterogeneity in the frequency of PQS in 11883 mtDNA**
 161 **sequences across the taxonomic sub-groups.** The frequency of PQS relative to genome
 162 length expressed as PQS/kbp. The organisms are grouped by Kingdom into the Protista
 163 (blue), Plantae (green), Fungi (yellow), and Animalia (red). The lines within the body of the
 164 violin plots represent the median and upper and lower quartiles. The range and distribution of
 165 sequences are also displayed. Genomes were obtained from the NCBI organelle genome
 166 database. Frequencies across the groups were found to be significantly different to one
 167 another ($P < 0.0001$) as determined via a non-parametric Kruskal-Wallis test.

168

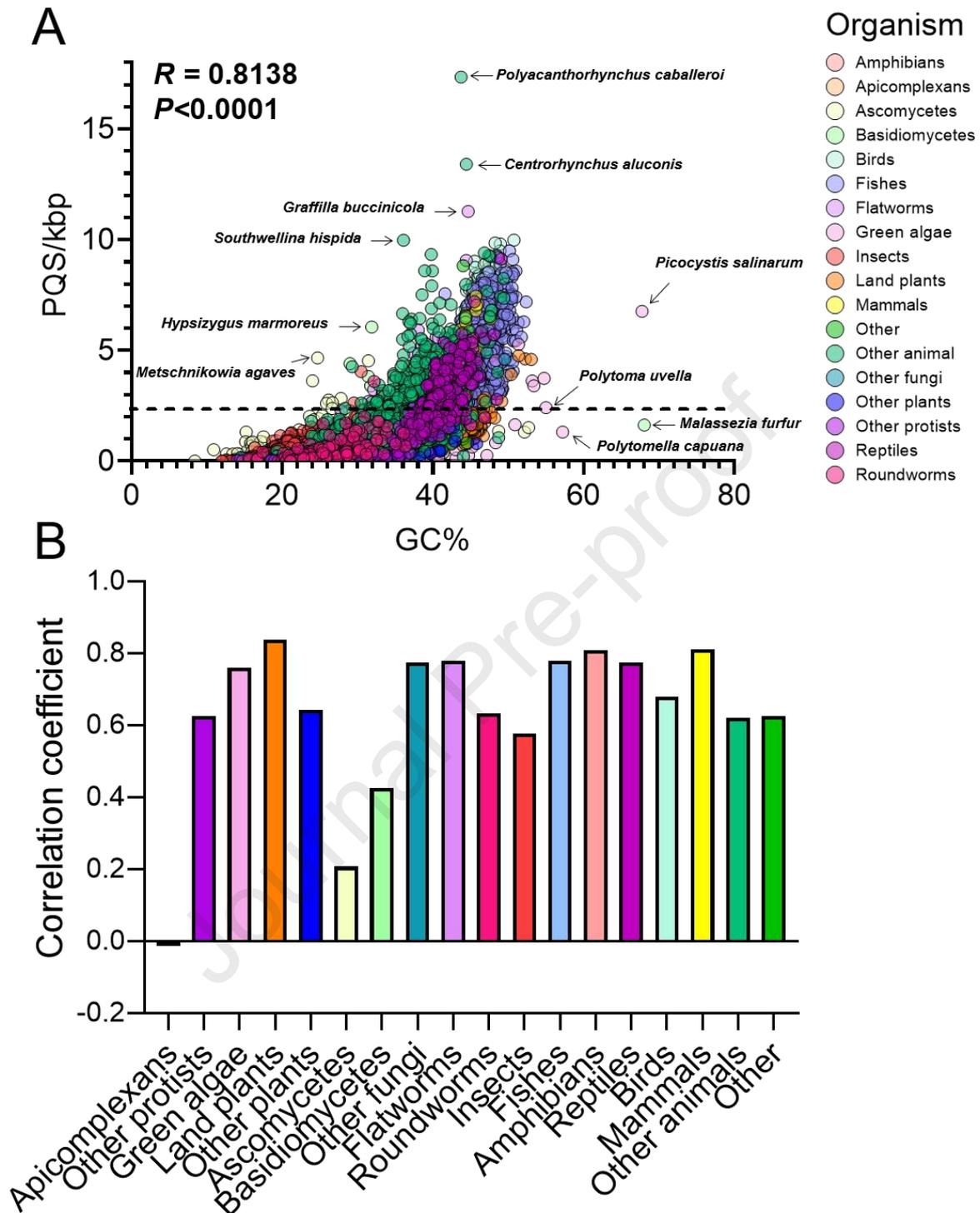
169 3.2 PQS frequency is correlated with mtDNA GC content

170 As quadruplexes can be found in GC-rich genomic regions, we next explored whether PQS
 171 frequency was correlated with genome GC content in mtDNA. In general, PQS frequency was
 172 found to be positively correlated with the genome GC content ($R=0.8138$; **Figure 2A**). This
 173 positive correlation was observed in all sub-groups except for the apicomplexans (**Figure 2B**;
 174 **Supplementary Figure S3**). However, some organisms had higher or lower PQS frequencies

175 relative to their GC content which deviated from the average. Some of these organisms with
176 high PQS frequencies included the parasites *Polyacanthorhynchus caballeroi*, *Centrorhynchus*
177 *aluconis*, and *Southwellina hispida*, the fungus *Metschnikowia agaves*, and flatworm *Graffilla*
178 *buccinicola* (**Figure 2A**). Those with lower PQS frequencies included the fungi *Malassezia*
179 *furfur*, *Candida gigantensis*, and *Candida subhashii*, and the green algae *Polytoma uvella*,
180 *Picocystis salinarum*, and *Polytomella capuana* (**Figure 2A**).

181

Journal Pre-proof



182

183 **Figure 2. The frequency of PQS is positively correlated with genome GC content. (A)**
 184 *There was a positive correlation between PQS frequency and genome GC content amongst*
 185 *all genomes used in the study. (B) The correlation coefficients for each group. Some*
 186 *organisms which had greater or fewer PQS than expected relative to their GC content are*
 187 *highlighted. The dashed line represents the average PQS/kbp for all organisms used in the*
 188 *study. Correlation was calculated via two-tailed Pearson's correlation coefficient ($P < 0.0001$).*

189

190 3.3 Quadruplexes are found in critical regulatory regions in mtDNA

191 To garner some insight into the potential roles that quadruplexes may play in mtDNA, we
192 endeavoured to identify where in the genome PQS were located. The genomic locations
193 discussed here were categorised based upon the annotation definitions used in the NCBI
194 database and we identified the PQS frequency in regions 100 bp before, within, and 100 bp
195 after the annotated feature.

196
197 Examining all the mitochondrial genomes together highlighted that PQS were found at greater
198 than average frequencies in key regulatory regions of mtDNA, including the 3'UTRs, D-loops,
199 replication origins, and stem loops (**Figure 3**). Conversely, PQS were found to be depleted in
200 the exons, introns, messenger RNA, and transfer RNA (**Figure 3**). However, when looking at
201 each subgroup individually, the picture is more complex and the distribution of PQS in genomic
202 features amongst the subgroups was highly varied (**Supplementary Table S2 and**
203 **Supplementary Figure S3**). Other regions of interest where PQS could be found at a high
204 frequency included the repeat regions (ascomycetes, basidiomycetes, flatworms, green algae,
205 other animals, other fungi), and non-coding RNA (ascomycetes, basidiomycetes, land plants;
206 **Supplementary Figure S3**).

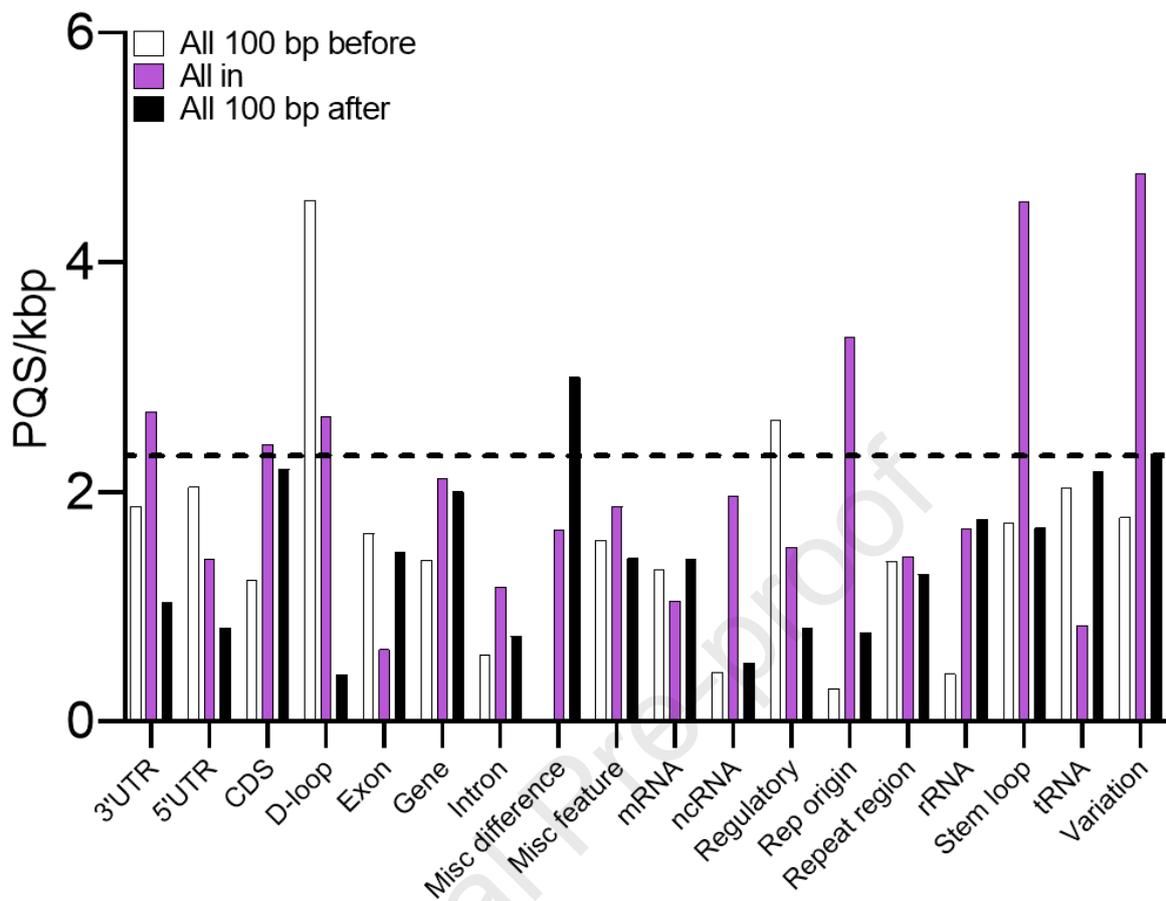
207
208 MtDNA G4s have been found with increased prevalence within cancer cells, to be associated
209 with DNA damage, and found in regions linked with the formation of deletion breakpoints in
210 patients with mitochondrial disorders (Falabella *et al.*, 2019; Butler *et al.*, 2020; Dahal *et al.*,
211 2021). However, growing evidence supports a role for G4s in biologically relevant
212 mitochondrial functions. Falabella and colleagues found that the G4-stabilising compound
213 3,11-Difluoro-6,8,13-trimethylquino[4,3,2-*k*]acridinium methylsulfate (RHPS4) could
214 preferentially localise to mitochondria in non-cancerous cells to bind to and stabilise
215 mitochondrial G4s (Falabella *et al.*, 2019). Depending upon the concentration used and the

216 amount of G4 stabilisation, this compound could either induce DNA damage, modulate
217 replication, or limit transcription. Our observations that PQS are found within regulatory
218 regions such as the replication origins and stem loops, in addition to their localisation
219 immediately before and within the genes themselves, is indicative of a potential regulatory role
220 within mtDNA and is supportive of these experimental observations. Notably, G4-stabilisation
221 by high concentrations of RHPS4 resulted in decreased strand-specific RNA abundance at
222 the D-loop (Falabella et al., 2019). The D-loop is the site of first strand replication and
223 stabilisation of G4s in this region could significantly modulate mitochondrial function. Thus, we
224 explored the mtDNA genomes further to identify potential G4s in the D-loops that might
225 indicate conserved mitochondrial functions in non-human organisms.

226

227

228



229

230 **Figure 3. PQS are highly represented in key regulatory regions of mtDNA in all the**
 231 **genomes analysed.** The frequency of PQS (PQS/kbp) within an annotated genomic feature,
 232 or within ± 100 bp of the feature was quantified. The dashed line represents the average
 233 PQS/kbp of all the genomes. Features which were found fewer than 10 times were omitted for
 234 clarity. Annotations were obtained from the NCBI database and data represents the average
 235 for all 11883 genomes analysed.

236

237 3.4 The D-loop/control region quadruplex sequence displays significant length 238 heterogeneity and variations are conserved in birds, fishes, reptiles, amphibians, and 239 mammals

240 Several quadruplexes have recently been identified to form in human mtDNA but the best
 241 studied quadruplex-forming sequence in human mtDNA (GCGGGGGAGGGGGGGTTTG)
 242 falls within the CSBII region of the D-loop [22,23; Figure 4A].

243

244 Although this sequence has previously been well characterised in human mtDNA and
245 identified in mice, similar sequences have not yet been fully explored in other organisms. In
246 support of its potential importance in mitochondrial biology, we found that this sequence (or
247 variations of this sequence) could be found highly conserved in D-loops and control regions
248 amongst the mammal (83.7%), amphibian (98.0%), bird (97.0%), reptile (95.8%), and fishes
249 (96.2%) mtDNA sequences analysed (**Supplementary Table S3**). This sequence was noted
250 if it fulfilled two criteria. First, the sequence had to be located within the D-loop region, and
251 second, it had to be formed of a contiguous run of guanines or two guanine tracts separated
252 by a short loop sequence composed of adenine and thymine. These sequences could not be
253 found in all organisms, but this may have been a limitation of genome quality, rather than that
254 these organisms lacked these sequences. However, equivalent sequences seemed to be
255 practically absent in all other subgroups.

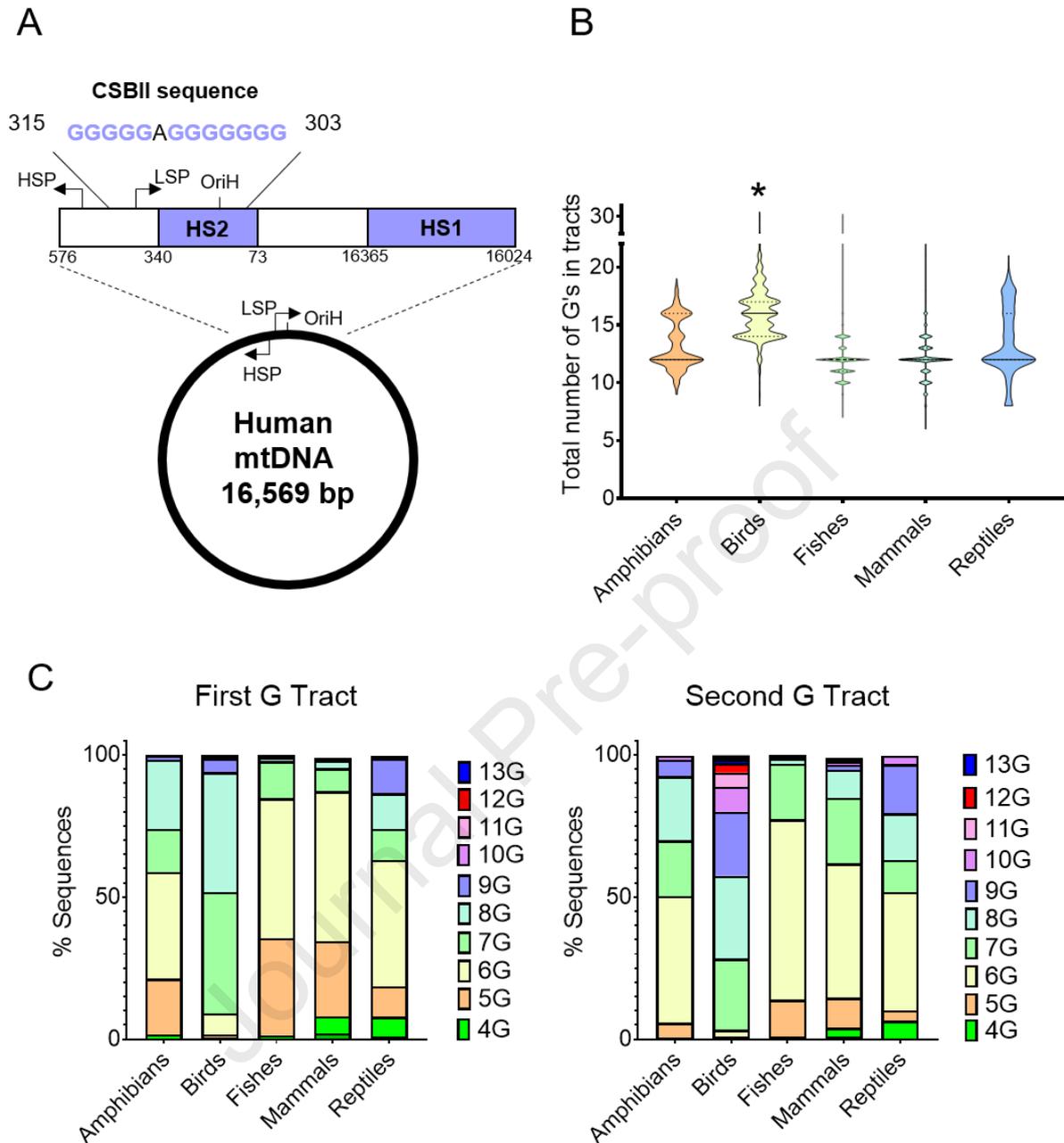
256

257 This G4-forming sequence was also found in the duplicated control regions of some birds and
258 reptiles, which suggests this G4 may have been evolutionary favoured and likely to play
259 important roles in mtDNA (**Supplementary Table S3**). Duplicated control regions, particularly
260 in parrots, have been shown to provide a selective advantage due to more efficient initiation
261 of replication or transcription and more replicating genomes per organelle ([24]. Consequently,
262 birds with replicated genomes live longer, have larger body masses, and are predisposed to
263 a more active flight. Again, when considering the potential detrimental effects that G4
264 stabilisation may have to mitochondria and that this region is prone to mutation, the exact
265 reason why sequences with potential to form these structures are so strongly conserved, or
266 even favoured, is unknown [25]. A potential explanation may be linked with the interaction
267 between the G4 in CSBII and RNA which forms an R-loop stabilising hybrid G4 [23]. The
268 mitochondrial R-loop plays critical roles in the replication, organisation, and expression of
269 mtDNA and compromising R-loop formation can result in mtDNA aggregation and disease
270 [6,26]. Therefore, G4s may also be a key participant in mtDNA and segregation through the

271 stabilisation of R-loops. The presence of G4s has been shown to result in premature
272 transcription of POLRMT but has been suggested to form R-loop structures which provide free
273 3' ends to prime subsequent DNA synthesis [27–29]. Recent evidence also indicates that the
274 G4-stabilised R-loop leads to increased transcription through a mechanism involving
275 successive rounds of R-loop formation [30].

276

277 Significant heterogeneity in the number of guanines in both the first and second G tracts of
278 this D-loop sequence were observed amongst all organisms and we noted at least 106
279 different G tract length combinations (**Supplementary Table S4**). In general, the second G
280 tract was found to be longer than the first, as seen previously, and the most common
281 combination observed throughout all organisms was 6 guanines in each tract. However, the
282 length of the first G-tract was not always shorter, as has been observed for humans [31].
283 Moreover, there was significant length heterogeneity in this region with G tract lengths ranging
284 from 2 to 21 and from 3 to 22 guanine residues in the first and second G tracts, respectively
285 (**Supplementary Table S4**). There were also large differences found in the linking sequence,
286 but the most frequently observed linking sequences between the G tracts were A, TA, or TTA
287 (**Supplementary Table S3**).



288

289 **Figure 4. A D-loop quadruplex sequence in humans is conserved in amphibians, birds,**
 290 **fishes, reptiles, and other mammals. (A)** The G4-forming sequence between nucleotides
 291 315 and 303 in the CSBII region of the D-loop in human mtDNA. The control region is indicated
 292 by dashed lines and the locations of the light and heavy strand promoters, the origin of H-
 293 strand replication, region of hypervariable sequence 1 and 2 (HS1 and HS2) are indicated.
 294 Adapted from Tan et al. 2006. **(B)** The total number of guanines in the first and second G-
 295 tracts combined for each group. The lines within the violin body represent the median (solid)
 296 and upper and lower quartiles (dashed). Birds had significantly more guanines in these tracts
 297 compared to the other groups ($P < 0.0001$) as determined by a Kruskal-Wallis test with Dunn's
 298 multiple comparisons. **(C)** The percentage of organisms in each group with first or second G-
 299 tract lengths of between 4 and 13 guanines. G-tract lengths outside of these ranges were less
 300 common and these were omitted to focus on the predominant groups.

301 Birds not only have the greatest frequency of PQS within their mtDNA, on average they also
302 have significantly more guanines in these D loop G tracts combined (**Figure 4B**). When
303 comparing the G tract lengths in detail, the most frequent G tract length observed in both the
304 first and second G tracts of mammals, fishes, amphibians, and reptiles was six (**Figure 4C**).
305 However, in birds, the most frequent G tract lengths were longer and were found to be 7 and
306 8 for the first and second G tracts, respectively (**Figure 4C**). Length heterogeneity of the G-
307 tracts in the D-loop/control region has been found to be associated with the amount of
308 transcription termination, with longer G-tracts resulting in increased termination [31].
309 Interestingly, increased length of these G-tracts is found favoured in cells with elevated growth
310 characteristics and it may be that increased transcription termination is associated with higher
311 levels of mtDNA replication [31].

312

313 **4. Conclusions**

314 Taken together, the conservation and prevalence of quadruplex-forming sequences in mtDNA,
315 and the D-loop is indicative of potential key regulatory roles for quadruplexes within mtDNA.
316 This study provides an in-depth overview of the similarities and differences between mtDNA
317 in highly diverse organisms and an insight into the importance of quadruplexes in their
318 genomes. Longer G-tract lengths and presence of G4-forming sequences in duplicated control
319 regions could be evolutionarily favoured in organisms which require higher levels of mtDNA
320 replication, such as birds which require increased replication of mitochondria to fulfil the
321 energy requirements necessary for flight, for example. However, we still have much to
322 discover, and the roles of quadruplexes in mtDNA require much greater examination and
323 experimental validation.

324

325 **Acknowledgements**

326 This work was supported by The Czech Science Foundation (18-15548S).

327 Conflict of interest

328 The authors disclose no conflicts of interest.

329 Author contributions

330 S.B. and V. conceived the study, N.B., M.D., and S.B. collected the data. N.B., M.D., and

331 S.B. analysed the data. N.B., V.B., and S.B. wrote the manuscript, V.B., acquired the

332 funding. All authors have approved the final article.

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349 **References**

- 350 [1] J.R. Friedman, J. Nunnari, Mitochondrial form and function., *Nature*. 505 (2014) 335–
351 343. <https://doi.org/10.1038/nature12985>.
- 352 [2] E. Murphy, H. Ardehali, R.S. Balaban, F. DiLisa, G.W. 2nd Dorn, R.N. Kitsis, K. Otsu,
353 P. Ping, R. Rizzuto, M.N. Sack, D. Wallace, R.J. Youle, Mitochondrial Function,
354 Biology, and Role in Disease: A Scientific Statement From the American Heart
355 Association., *Circ. Res.* 118 (2016) 1960–1991.
356 <https://doi.org/10.1161/RES.000000000000104>.
- 357 [3] G.S. Gorman, P.F. Chinnery, S. DiMauro, M. Hirano, Y. Koga, R. McFarland, A.
358 Suomalainen, D.R. Thorburn, M. Zeviani, D.M. Turnbull, Mitochondrial diseases., *Nat.*
359 *Rev. Dis. Prim.* 2 (2016) 16080. <https://doi.org/10.1038/nrdp.2016.80>.
- 360 [4] V. Brázda, R.C. Laister, E.B. Jagelská, C. Arrowsmith, Cruciform structures are a
361 common DNA feature important for regulating biological processes., *BMC Mol. Biol.*
362 12 (2011) 33. <https://doi.org/10.1186/1471-2199-12-33>.
- 363 [5] A. Herbert, Z-DNA and Z-RNA in human disease., *Commun. Biol.* 2 (2019) 7.
364 <https://doi.org/10.1038/s42003-018-0237-x>.
- 365 [6] J.M. Santos-Pereira, A. Aguilera, R loops: new modulators of genome dynamics and
366 function., *Nat. Rev. Genet.* 16 (2015) 583–597. <https://doi.org/10.1038/nrg3961>.
- 367 [7] J. Robinson, F. Raguseo, S.P. Nuccio, D. Liano, M. Di Antonio, DNA G-quadruplex
368 structures: more than simple roadblocks to transcription?, *Nucleic Acids Res.* (2021).
369 <https://doi.org/10.1093/nar/gkab609>.
- 370 [8] S. Burge, G.N. Parkinson, P. Hazel, A.K. Todd, S. Neidle, Quadruplex DNA:
371 sequence, topology and structure., *Nucleic Acids Res.* 34 (2006) 5402–5415.
372 <https://doi.org/10.1093/nar/gkl655>.
- 373 [9] S. Balasubramanian, L.H. Hurley, S. Neidle, Targeting G-quadruplexes in gene

- 374 promoters: a novel anticancer strategy?, *Nat. Rev. Drug Discov.* 10 (2011) 261–275.
375 <https://doi.org/10.1038/nrd3428>.
- 376 [10] Y. Zhang, M. Yang, S. Duncan, X. Yang, M.A.S. Abdelhamid, L. Huang, H. Zhang,
377 P.N. Benfey, Z.A.E. Waller, Y. Ding, G-quadruplex structures trigger RNA phase
378 separation., *Nucleic Acids Res.* 47 (2019) 11746–11754.
379 <https://doi.org/10.1093/nar/gkz978>.
- 380 [11] V. Brázda, J. Kolomazník, J. Lýsek, M. Bartas, M. Fojta, J. Šťastný, J.-L. Mergny,
381 G4Hunter web application: a web server for G-quadruplex prediction., *Bioinformatics.*
382 35 (2019) 3493–3495. <https://doi.org/10.1093/bioinformatics/btz087>.
- 383 [12] A. Bedrat, L. Lacroix, J.-L. Mergny, Re-evaluation of G-quadruplex propensity with
384 G4Hunter., *Nucleic Acids Res.* 44 (2016) 1746–1759.
385 <https://doi.org/10.1093/nar/gkw006>.
- 386 [13] R. Suzuki, H. Shimodaira, Pvclost: an R package for assessing the uncertainty in
387 hierarchical clustering., *Bioinformatics.* 22 (2006) 1540–1542.
388 <https://doi.org/10.1093/bioinformatics/btl117>.
- 389 [14] J. Cechová, J. Lýsek, M. Bartas, V. Brázda, Complex analyses of inverted repeats in
390 mitochondrial genomes revealed their importance and variability., *Bioinformatics.* 34
391 (2018) 1081–1085. <https://doi.org/10.1093/bioinformatics/btx729>.
- 392 [15] A.L. Mikheikin, A.Y. Lushnikov, Y.L. Lyubchenko, Effect of DNA supercoiling on the
393 geometry of holliday junctions., *Biochemistry.* 45 (2006) 12998–13006.
394 <https://doi.org/10.1021/bi061002k>.
- 395 [16] J.-N. Yang, A. Seluanov, V. Gorbunova, Mitochondrial inverted repeats strongly
396 correlate with lifespan: mtDNA inversions and aging., *PLoS One.* 8 (2013) e73318.
397 <https://doi.org/10.1371/journal.pone.0073318>.
- 398 [17] I. Skujina, R. McMahon, V.P. Lenis, G. V Gkoutos, M. Hegarty, Duplication of the

- 399 mitochondrial control region is associated with increased longevity in birds., *Aging*
400 (Albany, NY). 8 (2016) 1781–1789. <https://doi.org/10.18632/aging.101012>.
- 401 [18] M. Falabella, R.J. Fernandez, F.B. Johnson, B.A. Kaufman, Potential Roles for G-
402 Quadruplexes in Mitochondria., *Curr. Med. Chem.* 26 (2019) 2918–2932.
403 <https://doi.org/10.2174/0929867325666180228165527>.
- 404 [19] T.J. Butler, K.N. Estep, J.A. Sommers, R.W. Maul, A.Z. Moore, S. Bandinelli, F.
405 Cucca, M.A. Tuke, A.R. Wood, S.K. Bharti, D.F. Bogenhagen, E. Yakubovskaya, M.
406 Garcia-Diaz, T.A. Guillian, A.K. Byrd, K.D. Raney, A.J. Doherty, L. Ferrucci, D.
407 Schlessinger, J. Ding, R.M. Brosh, Mitochondrial genetic variation is enriched in G-
408 quadruplex regions that stall DNA synthesis in vitro., *Hum. Mol. Genet.* 29 (2020)
409 1292–1309. <https://doi.org/10.1093/hmg/ddaa043>.
- 410 [20] S. Dahal, H. Siddiqua, V.K. Katapadi, D. Iyer, S.C. Raghavan, Characterization of G4
411 DNA formation in mitochondrial DNA and their potential role in mitochondrial genome
412 instability., *FEBS J.* (2021). <https://doi.org/10.1111/febs.16113>.
- 413 [21] M. Falabella, J.E. Kolesar, C. Wallace, D. de Jesus, L. Sun, Y. V Taguchi, C. Wang,
414 T. Wang, I.M. Xiang, J.K. Alder, R. Maheshan, W. Horne, J. Turek-Herman, P.J.
415 Pagano, C.M. St Croix, N. Sondheimer, L.A. Yatsunyk, F.B. Johnson, B.A. Kaufman,
416 G-quadruplex dynamics contribute to regulation of mitochondrial gene expression.,
417 *Sci. Rep.* 9 (2019) 5605. <https://doi.org/10.1038/s41598-019-41464-y>.
- 418 [22] R.M. Andrews, I. Kubacka, P.F. Chinnery, R.N. Lightowlers, D.M. Turnbull, N. Howell,
419 Reanalysis and revision of the Cambridge reference sequence for human
420 mitochondrial DNA., *Nat. Genet.* 23 (1999) 147. <https://doi.org/10.1038/13779>.
- 421 [23] P.H. Wanrooij, J.P. Uhler, Y. Shi, F. Westerlund, M. Falkenberg, C.M. Gustafsson, A
422 hybrid G-quadruplex structure formed between RNA and DNA explains the
423 extraordinary stability of the mitochondrial R-loop., *Nucleic Acids Res.* 40 (2012)
424 10334–10344. <https://doi.org/10.1093/nar/gks802>.

- 425 [24] A.D. Urantówka, A. Krocak, T. Silva, R.Z. Padrón, N.F. Gallardo, J. Blanch, B.
426 Blanch, P. Mackiewicz, New Insight into Parrots' Mitogenomes Indicates That Their
427 Ancestor Contained a Duplicated Region., *Mol. Biol. Evol.* 35 (2018) 2989–3009.
428 <https://doi.org/10.1093/molbev/msy189>.
- 429 [25] M. Stoneking, Hypervariable sites in the mtDNA control region are mutational
430 hotspots., *Am. J. Hum. Genet.* 67 (2000) 1029–1032. <https://doi.org/10.1086/303092>.
- 431 [26] G. Akman, R. Desai, L.J. Bailey, T. Yasukawa, I. Dalla Rosa, R. Durigon, J.B.
432 Holmes, C.F. Moss, M. Mennuni, H. Houlden, R.J. Crouch, M.G. Hanna, R.D.S.
433 Pitceathly, A. Spinazzola, I.J. Holt, Pathological ribonuclease H1 causes R-loop
434 depletion and aberrant DNA segregation in mitochondria., *Proc. Natl. Acad. Sci. U. S.*
435 *A.* 113 (2016) E4276-85. <https://doi.org/10.1073/pnas.1600537113>.
- 436 [27] T.A. Brown, A.N. Tkachuk, D.A. Clayton, Native R-loops persist throughout the mouse
437 mitochondrial DNA genome., *J. Biol. Chem.* 283 (2008) 36743–36751.
438 <https://doi.org/10.1074/jbc.M806174200>.
- 439 [28] B. Xu, D.A. Clayton, RNA-DNA hybrid formation at the human mitochondrial heavy-
440 strand origin ceases at replication start sites: an implication for RNA-DNA hybrids
441 serving as primers., *EMBO J.* 15 (1996) 3135–3143.
- 442 [29] D. Kang, K. Miyako, Y. Kai, T. Irie, K. Takeshige, In vivo determination of replication
443 origins of human mitochondrial DNA by ligation-mediated polymerase chain reaction.,
444 *J. Biol. Chem.* 272 (1997) 15275–15279. <https://doi.org/10.1074/jbc.272.24.15275>.
- 445 [30] C.-Y. Lee, C. McNerney, K. Ma, W. Zhao, A. Wang, S. Myong, R-loop induced G-
446 quadruplex in non-template promotes transcription by successive R-loop formation.,
447 *Nat. Commun.* 11 (2020) 3392. <https://doi.org/10.1038/s41467-020-17176-7>.
- 448 [31] B.G. Tan, F.C. Wellesley, N.J. Savery, M.D. Szczelkun, Length heterogeneity at
449 conserved sequence block 2 in human mitochondrial DNA acts as a rheostat for RNA

450 polymerase POLRMT activity., Nucleic Acids Res. 44 (2016) 7817–7829.

451 <https://doi.org/10.1093/nar/gkw648>.

452

Journal Pre-proof

Highlights

- PQS frequency decreases with an increase in evolutionary distance
- PQS are over-represented in the 3'UTR, D-loops, replication origins, and stem loops
- Variation of G4 sequence in the D-loop is conserved across taxonomic sub-groups
- D-loop sequence is conserved in duplicated control regions of birds and reptiles
- Significant length heterogeneity in guanine tracts of the conserved D-loop sequence

Journal Pre-proof

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Journal Pre-proof