



HAL
open science

Comparative genomic analysis identifies small open reading frames (sORFs) with peptide-encoding features in avian 16S rDNA

Mathieu Mortz, Cyril Degletagne, Caroline Romestaing, Claude Duchamp

► **To cite this version:**

Mathieu Mortz, Cyril Degletagne, Caroline Romestaing, Claude Duchamp. Comparative genomic analysis identifies small open reading frames (sORFs) with peptide-encoding features in avian 16S rDNA. *Genomics*, 2020, 112 (2), pp.1120-1127. 10.1016/j.ygeno.2019.06.026 . hal-02345550

HAL Id: hal-02345550

<https://univ-lyon1.hal.science/hal-02345550v1>

Submitted on 13 Sep 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 Comparative genomic analysis identifies
2 Small Open Reading Frames (sORFs)
3 with peptide-encoding features in avian
4 16S rDNA

5

6 Mathieu Mortz, Cyril Dégletagne, Caroline Romestaing, Claude
7 Duchamp*

8

9

10 *Université de Lyon, Laboratoire d'Ecologie des Hydrosystèmes Naturels et Anthropisés, UMR 5023*
11 *CNRS, Université Claude Bernard Lyon 1, ENTPE, Villeurbanne cedex, France*

12

13 *Corresponding author: *Laboratoire d'Ecologie des Hydrosystèmes Naturels et Anthropisés, UMR*
14 *5023 CNRS, Université Claude Bernard Lyon 1, ENTPE, 43 Bvd 11 Novembre 1918, F-69622*
15 *Villeurbanne cedex, France.*

16 *Email: claude.duchamp@univ-lyon.fr*

17

18

19

20

21

22

23

24 **Short title:** Putatively functional sORFs in avian 16S rDNA

25

26

27

28

29

30

31

32 **Abstract**

33 The mitochondrial genome (mt-DNA) functional repertoire has recently been enriched in
34 mammals by the identification of functional small open reading frames (sORFs) embedded in
35 ribosomal DNAs. Through comparative genomic analyses the presence of putatively
36 functional sORFs was investigated in birds. Alignment of available avian mt-DNA sequences
37 revealed highly conserved regions containing four putative sORFs that presented low
38 insertion/deletion polymorphism rate (<0.1%) and preserved in frame start/stop codons in
39 >80% of species. Detected sORFs included avian homologs of human Humanin and Short-
40 Humanin-Like-Peptide 6 and two new sORFs not yet described in mammals. The amino-acid
41 sequences of the four putative encoded peptides were strongly conserved among birds, with
42 amino-acid p-distances (5.6 to 25.4%) similar to those calculated for typical avian mt-DNA-
43 encoded proteins (14.8%). Conservation resulted from either drastic conservation of the
44 nucleotide sequence or negative selection pressure. These data extend to birds the
45 possibility that mitochondrial rDNA may encode small bioactive peptides.

46

47

48 **Keywords:** Mitochondrial genome; Ribosomal RNAs; Humanin; Short-Humanin-Like
49 Peptides; SHLP6

50

51

52

53

54

55

56 **1. Introduction**

57 Mitochondria are cell organelles mainly involved in the production of ATP required to
58 sustain biological exergonic processes (1). ATP synthesis mostly occurs within the oxidative
59 phosphorylation (OXPHOS) system that is composed of five protein complexes combining
60 subunits encoded by nuclear (nuc-DNA) and mitochondrial (mt-DNA) genomes (2-3). In
61 animals, mt-DNA usually contains 37 genes that encode 22 tRNAs, 2 rRNAs (12S & 16S)
62 and 13 mitochondrial protein-coding genes (mt-PCGs) that complement the subunits
63 encoded by nuc-DNA to constitute the OXPHOS system (4).

64 This traditional mt-DNA description has been questioned in the last decade with the
65 description in various distantly related multicellular organisms of highly conserved small open
66 reading frames (sORFs), scattered within mt-DNA and showing characteristics suggesting
67 encoding function (5). In human mt-DNA, eight sORFs have already been identified as
68 effectively encoding small peptides and called Humanin (Hn) (6-7), MOTS-c (8) or SHLP1-6
69 (9). These peptides encoded by sORFs of different sizes (45-90 pb) embedded on both
70 strands of the 12S (MOTS-c) or 16S (Hn, SHLP1-6) rDNAs are obtained after mitochondrial
71 or cytosolic translation (10). Translated products are detected in the bloodstream and exert a
72 large spectrum of biological actions including anti-apoptotic and metabolic effects (11).
73 Owing to the large number of identified peptides and their biological roles, a novel concept of
74 mitochondria-derived peptide-signalling is slowly emerging in recent years (10, 12). To date
75 studies focused on human and rodent laboratory models, suggesting that this concept may
76 be extended at least to the class of mammals. However the existence of a dual encoding role
77 (rRNA and small peptides) for mitochondrial rDNAs has rarely been explored in other classes
78 of animals. Only one study, using data mining approaches, reported that the humanin gene
79 was conserved in 10 bird species while it was pseudogenized in many vertebrates (13). The
80 common presence of the humanin gene in primates, some rodents and birds (13) warrants
81 investigating more deeply the putative presence of mt-DNA-encoded peptides in birds.

82 Here, we propose an original *in silico* blind approach to test the peptide-encoding
83 potential of 12S and 16S rDNAs in the highly diversified class of birds through the detection
84 of sORFs. Birds together with mammals exhibit endothermy that requires a tachymetabolism
85 that might necessitate preserving metabolically active signals. Functional sORFs should a
86 priori be found in highly conserved regions of 12S/16S rDNAs and should also fulfil
87 characteristics of functional coding regions such as reading frame maintenance and amino-
88 acid sequence conservation. Our bioinformatics analysis therefore included stringent filters
89 (14) to avoid the detection of false-positive sORFs.

90

91 **2. Material and Methods**

92

93 *2.1. Collection of mt-DNA sequences and alignment*

94 All available (at the time of the study) complete avian mt-DNA sequences (RefSeq)
95 were extracted from GenBank (552 species). Sampled species represented all avian orders
96 in various proportions (42% Passeriformes, 11% Galliformes, 5% Anseriformes, 5%
97 Pelecaniformes, 5% Gruiformes, 5% Psittaciformes, 5% Charadriiformes, 3%
98 Collumbiformes, 3% Apodiformes, and 0.2-2% each for all other orders). Collected
99 sequences of avian mt-DNA were aligned using MUSCLE (15) through MEGA07 software
100 (16).

101 *2.2. Identification of highly conserved rDNA regions*

102 Highly conserved rDNA regions were identified by comparing region conservation to
103 the average conservation of the rDNA alignment. Conservation (C) corresponds to the
104 proportion of invariable sites ($C=1-(S/L)$), where S is the number of variable sites and L is the
105 net number of analysed positions. Conservation threshold was fixed to keep the rDNA
106 regions that are conserved at least 10% more than the rest of the rDNA gene (rDNA region

107 (C) > entire rDNA (C) + 0.1). The window width of rDNA region was determined as the
108 minimum width that allowed a region to be significantly conserved at $\alpha=0.05$ (under the
109 hypergeometric distribution). To select only highly conserved regions, we kept detected
110 regions with $\alpha < 10^{-6}$. Calculations (17) were made using DNAsp v5.10.01 (18).

111 2.3. *sORFs detection*

112 sORFs were detected by submitting the consensus sequences (threshold 25%,
113 considering gaps) of highly conserved regions to the NCBI bioinformatics tool ORFFinder.
114 Both heavy and light mt-DNA strands were screened and potential translation was addressed
115 using mitochondrial or standard genetic codes. To avoid the detection of an excessive
116 number of potential interesting sequences, the minimal length of detected sORFs was set at
117 45 nucleotides (<https://www.ncbi.nlm.nih.gov/orffinder/>).

118 2.4. *Peptide-encoding features and sORF selection process.*

119 A multiallelic estimation of the level of insertion/deletion (InDel) polymorphism
120 (considering all insertion/deletion events) was made for the concatenate of the thirteen mt-
121 PCGs (in the order: ATP6/8, COX1/2/3, CytB, ND1/2/3/4/4L/5/6), for the two 12S and 16S
122 rDNAs and for all highly conserved sORFs detected. This step was defined to select sORFs
123 with a level of InDel polymorphism close to or lower than that measured for all mt-PCGs with
124 the assumption that coding sequences should have low InDel polymorphism to maintain
125 reading frame and conserve the translated a-a sequence. We then verified, for each of the
126 552 species studied, the effective presence of in frame (compared with initially detected
127 consensus sORF) start codons (“ATG” for standard genetic code; “ATG” or “ATA” for
128 mitochondrial genetic code) and stop codons (“TAA”, “TAG” or “TGA” for standard genetic
129 code; “TAA”, “TAG”, “AGA” or “AGG” for mitochondrial genetic code). This enabled us to
130 estimate a percentage of sORF preservation that corresponds to the proportion of species
131 with a given putative encoding sORFs ≥ 45 nucleotides flanked by in frame start/stop codons.
132 A high degree of preservation is expected for functional and biologically relevant peptide

133 encoding sORFs. A stringent threshold of $\geq 80\%$ of preservation among species led to select
134 only a limited number of sORFs that were compared with human homologous mt-DNA
135 regions. Comparison was performed by aligning avian 16S consensus rDNA with human 16S
136 rDNA reference sequence (Homo sapiens mitochondrion complete genome, NCBI Reference
137 Sequence NC_012920.1, Position: 1671:3229). We also estimated for each of the remaining
138 selected sORFs the non-synonymous substitution rate (number of non-synonymous
139 substitutions / non-synonymous sites, K_a), the synonymous substitution rate (number of
140 synonymous substitutions / synonymous sites, K_s) (19-25) and both nucleotide and amino-
141 acid p-distances. Values were compared with those calculated for avian mt-PCGs. All
142 calculations were made through MEGA07 (16) and DNAsp v5.10.01 (18).

143 2.5. RT-PCR analysis of sORF expression.

144 sORF functionality was assessed by the detection of targeted transcripts using
145 reverse transcription-polymerase chain reaction (RT-PCR). Total RNA was extracted with
146 Trizol from either pectoralis muscle sample obtained from a 6-wk old Muscovy duckling
147 (*Cairina moschata*). After digestion of RNA extracts with DNase I (Thermo Fisher Scientific),
148 mRNA were purified using a Dynabeads® mRNA purification kit (Thermo Fisher Scientific)
149 and quantified using a Qubit™ Fluorometer (Thermo Fisher Scientific). Reverse transcription
150 (RT) of 1 μg RNA was performed using polyT as primer and 200 U MMLV-RT (Invitrogen)
151 according to the manufacturer protocol. PCR was then performed in a total volume of 20 μL
152 using platinum Taq polymerase (Invitrogen) and a Hybaid thermocycler (Thermo Fisher
153 Scientific). After an initial denaturation at 94°C for 3 minutes, 30 cycles of denaturation at
154 94°C for 45 s, annealing at 55°C for 30 s and extension at 72°C for 60 s were used. A final
155 extension was performed at 72°C for 10 min. Forward and reverse primers were designed
156 from muscovy duck 16S rDNA sequence (GenBank Accession: NC_010965.1) and were: av-
157 Hn 5'- CCT AAC CGT GCA AAG GTA GC-3' and 5'- AGT TCC ACA GGG TCT TCT CG -3'
158 (product size 155 bp); av-scO2 5'- ATA GAA ACC GAC CTG GAT TG -3' and 5'- GTT TAC
159 GAC CTC GAT GTT GG -3' (product size 127 bp). Actin was used as a positive control for

160 mRNA extracted from muscle 5'-GAC GAG GCC CAG AGC AAG AGA-3' and 5'-GGG TGT
161 TGA AGG TCT CAA ACA-3' (product size 225 bp). A positive control of mt-DNA transcription
162 was also used: av-Cox1 5'-AAC TAG GCC AAC CAG GGA CT-3' and 5'-AGA TGA GGC
163 GAG CAG TAG GA-3' (product size 230 bp). The amplified products were separated
164 according to their size on 2% agarose gel stained with ethidium bromide and images were
165 acquired with a Kodak Digital Science™ DC120 Camera.

166

167 **3. Results and Discussion**

168

169 *3.1. Identification of highly conserved regions & detection of avian sORFs*

170 The conservation pattern of avian rDNAs was assessed by calculating nucleotide
171 divergences for all positions without insertion/deletion event along the 12S and 16S rDNAs.
172 As shown in Fig.1, the average nucleotide divergence among birds of a 45 nucleotides
173 sliding window was heterogeneously distributed along 12S and 16S rDNAs. In both 12S
174 (Fig.1A) and 16S rDNAs (Fig.1B), divergence was heterogeneously distributed with regions
175 of relatively high divergence (>15%) and others with very low divergence (<5%), suggesting
176 that some regions are more functionally constrained among bird species. Heterogeneous
177 conservation along 12S/16S rDNA with strongly conserved positions has already been
178 described in several prokaryotic and eukaryotic taxa (27, 28), but present data are the first to
179 focus specifically on birds.

180 Heterogeneous conservation of rDNA sequence is not surprising as evolution of rDNA
181 depends on the secondary structure of encoded rRNA, with fast- and slow-evolving positions
182 in loop-encoding and stem-encoding regions respectively (29). Specific peptide-encoding
183 sORF in slow-evolving regions could constitute an additional constraint considerably
184 affecting the conservation of some specific portions of rDNA. Accordingly, our analysis on
185 bird sequences revealed three hotspots of remarkably small divergences in either 12S ($2.1 \pm$

186 0.9% at Pos. 508 (window 485-529)) or 16S rDNA ($1.4 \pm 1.0\%$ at Pos. 840 (window 819-863)
187 or $0.08 \pm 0.03\%$ Pos. 1303, (window 1282-1326)) (Fig.1). Assuming that important peptide
188 encoding rDNA regions should be strongly conserved during bird evolution, we looked for
189 highly significantly ($p < 10^{-6}$) conserved regions in 12S/16S rDNAs. Calculations (see
190 Methods) led to identify 5 significantly ($p < 0.05$) conserved regions and 3 of them, located
191 either in 12S or 16S rDNAs, were highly significantly ($p < 10^{-6}$) conserved among bird species
192 and logically contained the three hotspots of drastic conservation (Fig.1).

193 Using the avian consensus sequences of the identified highly conserved regions,
194 sORFs were then detected with ORFFinder. We searched for potential coding sequences
195 ≥ 15 a-a (≥ 45 nucleotides) with reference to vertebrate mitochondrial genetic code and/or
196 standard genetic code, according to the features noticed in mammals for the peptide-
197 encoding sORFs already identified (5, 6-9, 14). This strategy enabled us to detect 17 sORFs
198 embedded in either 12S (5 sORFs, Fig.1A) (Table S1) or 16S (12 sORFs, Fig.1B) (Table
199 S2). sORFs were numbered according to their position in the sequences, 1 to 5 in 12S and 6
200 to 17 in 16S and labelled as being detected by referring to either mitochondrial genetic code
201 (e.g. 1m/2m...) or standard genetic code (e.g. 1s/7s...). Two sORFs were detected by
202 referring to both mitochondrial and standard genetic codes (1m/1s/10m/10s) and were each
203 considered as distinct sORFs (mitochondrial vs cytosolic forms) for the following analyses,
204 because of genetic code-dependent variations in i) sORF length related to differences in
205 sORF initiation and stop codons positions and ii) in the composition of translated a-a
206 sequences.

207 3.2. *Low level of InDel polymorphism in avian sORFs?*

208 The mere fact of detecting sORFs in conserved regions does not mean that they are
209 functional and specifically conserved during evolution (30, 31), especially since avian rDNA
210 sequences contain a high percentage of GC both in 12S (49.5%) and 16S (45.1%) that
211 would likely favour the random presence of sORFs in non-coding sequences. However, a
212 direct consequence of the existence of functional peptide encoding regions is a weak level of

213 insertion / deletion (InDel) events to maintain reading frame and thus to conserve the
214 translated a-a sequence. This feature can be verified for avian mt-DNA genes by measuring
215 the InDel polymorphism rate (InDel Pi) for rDNAs and for the concatenate of all mt-PCGs.
216 The InDel Pi (/100 sites) was more than 20 times lower for mt-PCGs (0.16%) than that
217 measured for 12S (3.68%) or 16S (4.6%) rDNAs (Fig.2).

218 Within the detected conserved sORFs, nine sORFs presented InDel Pi in the same
219 order of magnitude (ranging from 0.96 to 3.55%) as those measured for the entire 12S or
220 16S rDNAs suggesting the absence of specific pressure during evolution to maintain the
221 reading frame by comparison with the rDNA gene in which they are included. By contrast,
222 seven sORFs (2_m, 9_m, 10_{ms}, 11_s, 15_m, 16_s, 17_m) presented very low InDel Pi (ranging from
223 0.007 to 0.094%) that were even lower than that measured for the concatenate of avian mt-
224 PCGs (Fig.2). These extremely low InDel Pi are consistent with the maintenance of putative
225 reading frames, as would be expected in peptide/protein-coding regions. Present data
226 therefore suggest that the seven sORFs detected in the consensus avian 16S rDNA might
227 encode peptides that impose functional constraints on sequences during evolution by
228 comparison with the rest of the gene in which they are included.

229 3.3. *In frame sORF preservation among bird species?*

230 To provide additional arguments in favour of their putative peptide-encoding role, we
231 focused on these seven detected sORFs and verified their effective presence in all 552 bird
232 species. For each sORF in every species, we looked for the presence of both start and stop
233 codons in frame with the initially detected consensus sORF and flanking a region of ≥ 45
234 nucleotides. We then calculated the percentage of sORF preservation (Fig.3) that
235 corresponds to the proportion of bird species meeting those criteria and kept sORFs
236 preserved in more than 80% of studied bird species. This high level of preservation was
237 chosen to keep sORFs with a significant biological role potentially imposing strong
238 conservation with evolution.

239 As shown in Fig.3, estimations of sORF preservation varied among the seven selected
240 sORFs indicating that neither the low nucleotide diversity (Fig.1) nor the low InDel
241 polymorphism rate (Fig.2) of the sORFs in the consensus sequence would ensure their
242 preservation in all bird species. For the last 12S sORF detected in the consensus sequence
243 and that fulfilled previous filters (2_m), only 41.3% of bird species preserved a ≥ 45 nucleotides
244 sORF. Similarly, despite strong nucleotide sequence conservation, sORF 9_m and 15_m were
245 preserved in 66.3% and 29.9% of birds, respectively. Only a limited number of sORFs were
246 preserved in more than 80% of studied species and they were all found in 16S rDNA: 10_m
247 (82.1%), 10_s (81.3%), 11_s (91.3%), 16_s (97.8%) and 17_m (99.1%). Such high degree of
248 preservation is indicative of functional constraints to maintain the sequence of specific
249 sORFs in most bird species. The presence of sORFs with robust peptide-encoding features
250 exclusively in avian 16S rDNA is consistent with the observations in mammals in which
251 identified encoding sORFs were mostly found within 16S rDNA ($n=7$) (6, 7, 9) while only one
252 encoding sORF was detected in 12S (8).

253 *3.4. Are avian sORFs potentially coding for orthologs of peptides already identified in* 254 *humans?*

255 The consensus sequences of the four remaining avian 16S sORFs were aligned with
256 human 16S rDNA sequence (NC_012920.1, Pos.1671:3229). Both avian consensus and
257 human sORF homologous regions were translated into a-a sequences using appropriate
258 genetic codes and compared (Table 1). Among the four remaining sORFs, we identified
259 putative avian orthologs of humanin (sORF 10_{ms} : av-Hn) and SHLP6 ((sORF 16_s : av-SHLP6)
260 that were encoded by the strand (+) of avian 16S rDNA. The deduced sequence of av-Hn
261 consensus peptide (obtained after either mitochondrial or cytosolic translation) showed 58%
262 identity (14 of the 24 a-a were identical) with the cytosolic form of human Hn. Present data
263 confirm and extend the previous finding of a potentially functional humanin gene in ten bird
264 species (13) by documenting the conservation of av-Hn in $>80\%$ of the 552 bird species in
265 which mt-DNA sequence was available. There was an even greater degree of a-a

266 conservation through evolution for SHLP6 as the human and avian proteins showed 90%
267 identity (18/20 identical a-a) (Table 1).

268 The presence of conserved orthologs for humanin and SHLP6 in birds and mammals
269 suggests that these peptides exhibit essential biological roles. In mammals, humanin was
270 shown to reduce apoptosis (32, 33) and to favour cardioprotection (34), while the various
271 SHLPs were shown to act as age-dependent regulators of apoptosis, insulin sensitivity and
272 inflammatory markers (9). The essential actions of these peptides as potent cell-survival and
273 metabolic factors might contribute, at least in part, to the necessity to preserve their a-a
274 composition and functional a-a motifs (Table 1) during evolution. It remains to be shown that
275 these biological roles have been preserved in birds.

276 The other two sORFs (11_s, 17_m) identified in our analysis were both putatively encoded
277 by the strand (-) of avian 16S rDNA and did not match with any characterized human
278 encoding sORFs. These new identified sORFs were called av-scO1 (11_s) and av-scO2 (17_m)
279 for “avian strongly conserved sORFs 1 and 2”. The presence of putative encoding sORFs in
280 the non-coding strand (-) of avian 16S rDNA is in agreement with the presence of sORFs
281 encoding SHLPs (SHLP1 to 5) in the non-coding strand (-) of human 16S rDNA (9).
282 Interestingly, despite a global overlap with av-Hn (pos. 886-960, strand (+)) and thus a same
283 substitution rate in their common nucleotide area (pos. 913-960), no homologous sORF was
284 identified in human for av-scO1. The translated human a-a sequence would result in a a-a
285 sequence with only 24% identity (4/17 a-a identical) with the avian putative peptide. For av-
286 scO2, the potential human peptide despite being shorter would show 94% identity with the
287 avian form (16/17 a-a identical). These observations suggest that av-scO1 appears to be
288 specific to bird species while av-scO2 would be a more ubiquitous sORF highly conserved at
289 least in endothermic vertebrates. The experimental confirmation that these sORFs are
290 functional and encode for biologically active peptides now remains to be appropriately
291 verified. To provide additional arguments concerning the potential encoding function of
292 detected sORFs we analysed the nucleotide composition around the start codon to detect

293 potential Kozak sequence (GCC(G/A)CCAUGG) that represent translation initiation signals
294 (TIS) facilitating the recognition of AUG codon by ribosomal subunits (35). No strong robust
295 Kozak-like sequences could be found in the detected sORFs. When focusing on the two
296 main sites (3(G/A) and +4G) associated with high translation efficiency in mammals (36), we
297 only noted the presence of +4G for av-Hn and av-scO2. Nevertheless, the absence of
298 characteristic TIS does not mean that selected sORFs could not encode peptides on account
299 of the fact that only 40% of mammalian mRNAs present highly efficient TIS and 10% of them
300 do not present the characteristic nucleotides -3(G/A) and +4G (37).

301 Since our large scale blind analysis successfully identified putative avian orthologs of
302 peptides already identified in mammals, it gave some confidence to tackle the
303 characterization of the peptides encoded by the new sORFs. To strengthen our analyses, we
304 harnessed available bird transcriptomic data (NCBI GEO DataSets) to determine whether the
305 putative sORF are transcribed. Most interestingly, we found sequences corresponding to the
306 four selected sORFs in mRNA libraries obtained from several tissues of *Gallus gallus* (38) or
307 *Anser cygnoides* (39) including heart, breast muscle or adipose tissue. Note that mRNA
308 libraries were prepared for sequencing using TruSeq stranded mRNA sample preparation
309 guide or the Ribo-Zero™ kit to remove rRNA. Although the possibility that transcripts
310 containing the sORFs encoded by the strand(+) could be related to traces of 16S rRNA not
311 removed by mRNA purification cannot be fully excluded, the finding of transcripts encoded by
312 the stand(-) such as av-scO1 and 2 is not expected unless they relate to specific
313 transcriptions of these regions (9). To further strengthen our in silico findings we performed
314 preliminary investigations of potential transcription using RT-PCR detection of one sORF of
315 each strand (av-Hn and av-scO1). As shown in Fig.4, RT-PCR amplicons of expected size
316 could be found using purified mRNAs extracted from duck pectoralis muscle. When RT was
317 omitted, no amplicon was detected showing the specific amplification of poly-adenylated
318 transcripts. Although these preliminary observations are promising signatures of sORF
319 functionality, the presence of translated peptides now requires to be investigated.

320 3.5. *Negative selection pressure for avian sORFs?*

321 The amino acid composition conservation of a peptide / protein coding gene among
322 species is known to result from a strong global nucleotide conservation, but also from an
323 increase of the inconsequential rate of synonymous (syn.) substitutions (number of syn.
324 substitution / number of syn. sites, K_s) at the expense of non-syn. substitutions rate (number
325 of non-syn. substitution / number of non-syn. sites, K_a). Thus, negative selection pressure
326 (K_a/K_s ratio <1) highlighted when comparing nucleotide sequence of different species is
327 generally the consequence of the presence of common peptide / protein coding region, which
328 requires to conserve the amino acid composition during evolution (40). To assess the type of
329 selection pressure underwent by the four selected 16S sORFs within bird species, both K_a
330 and K_s were estimated for each pairwise comparison of species and compared with that
331 underwent by the 13 mt-PCGs. K_a values were plotted as a function of K_s values for each
332 sORF and for the concatenate of avian mt-PCGs. With 552 bird species to compare, this
333 represented 152 352 pairwise comparisons per plot. To simplify the representation, only the
334 regression lines of each plot were drawn (Fig.5) and their respective slope (K_a/K_s) and
335 regression coefficient (R^2) were given. Results indicated that except for av-scO2 that
336 appeared to be under neutral selection (K_a/K_s slope = 0.94), av-scO1 (K_a/K_s slope = 0.32)
337 and the putative orthologs of humanin (K_a/K_s slope = 0.25 (mt.) / 0.22 (std.)) and SHLP6
338 (K_a/K_s slope = 0.36) were all under a strong negative selection pressure during bird
339 evolution. Such negative selection pressure is indicative of the need to conserve the amino
340 acid composition of the putative encoded peptides.

341 As shown in Fig.5, the concatenate of proteins encoded by avian mt-PCGs appeared
342 under an even higher (K_a/K_s slope = 0.10) negative selection than the selected avian sORFs
343 suggesting that the conservation of mt-proteins among bird species would be accordingly
344 higher. To verify this point, we calculated the amino acid conservation (p-distances) among
345 birds for the peptides/proteins encoded by the four sORFs and mt-PCGs. As shown in Fig.
346 5A, similar or greater conservations of a-a composition were unexpectedly found for the

347 translated peptide sequences of av-Hn ($15.4 \pm 4.8\%$ (mt.); $16.9 \pm 4.6\%$ (std.)), av-SHLP6
348 ($6.6 \pm 2.0\%$), av-scO1 ($25.4 \pm 4.9\%$) and av-scO2 ($5.6 \pm 2.0\%$) as compared with the
349 concatenate of all mt-PCGs ($14.8 \pm 0.3\%$) (Fig.6A). Despite overall conservation of
350 sequences, there were small substitutions in amino-acid composition that could affect the
351 functionality of the encoded peptide depending on the region concerned. It is difficult at this
352 stage to conclude about their impact on the functionality of the encoded peptides without
353 molecular validation of peptide function at the amino-acid level. However, the conservation of
354 the two putatively functional motifs found in av-Hn ("SEIDL") or av-scO2 ("TDLD") was
355 analysed and compared with the global p-distance measured for the entire sORF. We found
356 a remarkable conservation of the functional motif of av-scO2 TDLD ($1.0 \pm 0.2\%$) that
357 underwent a stronger constraint than the entire av-scO2 sORF ($5.6 \pm 2.0\%$), suggesting a
358 putative important role in the peptide functionality. By contrast, the av-Hn functional motif
359 SEIDL was not more conserved ($21.7 \pm 0.17\%$ (both mt. and std.)) than the entire sORF
360 ($15.4 \pm 4.8\%$ (mt.); $16.9 \pm 4.6\%$ (std.)). This was related to a strong positive selection
361 underwent by the first a-a of the SEIDL motif (mean Ks: 3.7%; mean Ka: 25.6%), while the
362 last four positions were less variable ($10.7 \pm 4.7\%$). Thorough analysis revealed the
363 predominance of only two amino-acids at the first position (the fourteenth position when
364 considering entire av-Hn peptide sequence) of the motif, with 291 species (53%) presenting
365 a serine (Ser14) and 228 species (41%) presenting a glycine (Gly14) at this position.
366 Interestingly, a Gly14-Hn variant was also characterized in mammals, and was shown to
367 have a 1000-fold higher neuroprotective activity than the Ser14-Hn variant (41). The non-
368 synonymous mutations observed in birds at the 14th position of av-Hn might thus reflect an
369 adaptation of av-Hn functionality in birds that encode a Gly14-Hn but additional experimental
370 verifications of these aspects are needed. Another functional site described in the literature
371 for Hn was the cysteine at the eighth position of the peptide (Cys8), observed in both avian
372 consensus and human peptide sequences. The presence of a cysteine at this position was
373 essential to the anti-apoptotic function of the peptide (42). This amino acid was found to be
374 present in nearly all 552 bird species studied, except one (*Turdus philomelos*) that presented

375 a frameshift that modified the encoding amino-acid at this position. More investigations are
376 needed to verify if this InDel event indeed occurred or if it is only the consequence of a
377 sequencing error.

378 Interestingly, the four putatively encoded peptides showed greater amino-acid
379 conservation than those measured for known mt-PCGs in contrast with calculated Ka/Ks
380 values (Fig.5). However, negative pressure intensity not only depends on the functional
381 importance of the sequence of the encoded protein in studied species but also on the global
382 substitution rate observed between nucleotide sequences (43-45). We therefore plotted
383 Ka/Ks ratio as a function of nucleotide conservation for the known or putative coding
384 sequences in birds (Fig.6B). There was a global negative relation between negative pressure
385 intensity and the nucleotide substitution rate estimated by measuring the nucleotide p-
386 distance for each studied regions (Fig.6B). For similar nucleotide conservation, for instance
387 with av-SHLP6 and av-scO2, Ka/Ks values could be very different. This might be due to the
388 fact that an increase of the global substitution rate in protein coding regions will mainly affect
389 synonymous sites to conserve a-a composition. Consequently, fast-evolving coding genes
390 show higher Ks values and consequently lower Ka/Ks ratio than slow-evolving genes,
391 regardless of the functional importance of the encoded a-a (43). Therefore, av-scO2,
392 drastically conserved among bird species (nuc. p-distance= $2.0\pm 0.9\%$) and resulting in a
393 strong peptide conservation ($5.6\pm 2.0\%$), could be qualified as a slow-evolving coding gene
394 as compared with the fast-evolving mt-PCGs (nuc. p-distance= $20.8\pm 0.2\%$) showing higher
395 protein divergence ($14.8\pm 0.3\%$). Thus, the apparent neutral selection observed for av-scO2
396 (Ka/Ks= 0.94) in Fig.5 was the consequence of extremely low values of both Ka (mean Ka=
397 3.2%) and Ks (mean Ks= 2.3%), as compared with the values calculated for the concatenate
398 of all mt-PCGs (mean Ka= 8.3% ; mean Ks= 55%). It was not related to the importance of av-
399 scO2 a-a composition. Altogether, present results indicate that the strong conservation of the
400 deduced a-a sequence for the four selected 16S sORFs, that was similar to that calculated

401 for avian mt-PCGs, was the result of drastic nucleotide conservation (av-SHLP6 and av-
402 scO2) or strong negative selection pressure (av-Hn and av-scO1) during bird evolution.

403 3.6. Are there only four putatively encoding sORFs in avian 12S and 16S rDNA?

404 Our cautious approach to detect sORFs verifying peptide-encoding features after
405 applying stringent filters aimed at minimizing the risk of selecting false-positive sequences.
406 The down side of this drastic selection is that some of the rejected sORFs, not fulfilling all
407 criteria, could anyway encode peptides. Filtering on InDel polymorphism was relevant
408 because it identified a distinct sORF population presenting very low InDel rates, in
409 agreement with the need to maintain reading frames specific to peptide/protein encoding
410 regions. Filtering on the presence of both start and stop codons in frame with initially
411 detected consensus sequences should however be moderated. Indeed, the conservative
412 vision which attributes to “AUG” (and “AUA” in mitochondria) the exclusive initiation of
413 translation is debated (46). As an example, when regarding all the avian mt-PCGs collected
414 for this study, we noticed a non-traditional initiation codon (other than “AUG” or “AUA”) in
415 15.5% of bird sequences. In a similar way, the condition of in frame start codon is not verified
416 for SHLP1 and SHLP3 in the 16S rDNA of Sprague Dawley rats (Table S3), while the
417 peptides are detected in blood and SHLP3 exhibits biological effects (9). We therefore
418 analyzed the selection pressure underwent by the three sORFs discarded due to their low
419 percentage of “sORF preservation” (2m, 9m, 15m). Interestingly, we found similar Ka/Ks
420 values to those found for the four selected sORFs that fulfilled all criteria, with a strong
421 negative selection for 12S-sORF 2m (Ka/Ks slope: 0.30) and 16S-sORF 9m (Ka/Ks slope:
422 0.15) allowing a good conservation of amino-acid composition ($7.8 \pm 4.7\%$ and $8.5 \pm 3.8\%$,
423 respectively) (Table S1/S2). The 16S-sORF 15m underwent a less potent negative selection
424 (Ka/Ks slope: 0.65), but the strong conservation of its nucleotide composition ($7.9 \pm 2.3\%$)
425 enabled a good conservation of amino-acid composition ($11.8 \pm 4.6\%$) (Table S2), as was
426 the case for av-SHLP6 and av-scO2. These 3 sORFs, even if they are not as robust
427 candidates as the four sORFs selected by our stringent cut-off selection, might also encode

428 peptide in bird species. It follows that additional molecular and biological approaches are
429 requested to overcome the limits of bioinformatics analysis (47) and draw up an exhaustive
430 list of all the peptides effectively encoded by mitochondrial rDNAs.

431 **4. Conclusion**

432 A large scale comparative genomic analysis of available avian mt-DNA sequences
433 (552 species) showed that in highly conserved regions of avian 16S rDNA, four sORFs
434 exhibited peptide-encoding features comparable to those measured for all avian mt-PCGs,
435 with a low InDel polymorphism rate (<0.1%), a high preservation of both start and stop
436 codons (present in more than 80% of species) and a high conservation of a-a composition
437 allowed by very low substitution rate and/or negative selection pressure on nucleotide
438 sequences. Among the four selected sORFs, two were identified as encoding for the putative
439 orthologs of human humanin and SHLP6 while the others were highly conserved at least
440 among birds suggesting important biological roles for the putative encoded peptides.
441 Altogether, these findings extend to birds the notion already characterised in human and
442 rodents that sORFs embedded in mitochondrial rDNA may encode biologically active
443 peptides.

444

445 **Author contributions**

446 MM collected data, aligned sequences, analyzed and interpreted data and prepared
447 the first draft of the manuscript. CDu and CR conceived the study and supervised the
448 analysis. CDe contributed to the analysis of data. All authors discussed the results and
449 implications, edited the manuscript and gave final approval for publication.

450

451 **Acknowledgements**

452 The authors thank Christophe Douady, Tristan Lefebure and Laurent Duret for helpful
453 discussions and suggestions. The project was financially supported by the Centre National

454 de la Recherche Scientifique and the University of Lyon. MM was in receipt of a fellowship
455 from the Ministère de l'Enseignement Supérieur et de la Recherche.

456

457 **Conflict of interest**

458 The authors report no conflict of interest.

459

460 **References**

- 461 [1] Papa S, Martino PL, Capitanio G, Gaballo A, De Rasmio D, Signorile A, Petruzzella V.
462 The oxidative phosphorylation system in mammalian mitochondria. *Adv Exp Med* 942:
463 3-37, 2012. doi:10.1007/978-94-007-2869-1_1.
- 464 [2] Wallace DC, Fan W, Procaccio V. Mitochondrial energetics and therapeutics. *Annu*
465 *Rev Pathol* 5: 297-348, 2010. doi:10.1146/annurev.pathol.4.110807.092314.
- 466 [3] Shokolenko IN, Alexeyev MF. Mitochondrial DNA: A disposable genome? *Biochim*
467 *Biophys Acta* 1852: 1805-1809, 2015. doi:10.1016/j.bbadis.2015.05.016.
- 468 [4] Boore JL. Animal mitochondrial genomes. *Nucleic Acids Res* 27: 1767-1780, 1999.
469 doi:10.1093/nar/27.8.1767.
- 470 [5] Breton S, Milani L, Ghiselli F, Guerra D, Stewart DT, Passamonti M. A resourceful
471 genome: updating the functional repertoire and evolutionary role of animal
472 mitochondrial DNAs. *Trends Genet* 30: 555-564, 2014. doi:10.1016/j.tig.2014.09.002.
- 473 [6] Hashimoto Y, Niikura T, Tajima H, Yasukawa T, Sudo H, Ito Y, Kita Y, Kawasumi M,
474 Kouyama K, Doyu M, Sobue G, Koide T, Tsuji S, Lang J, Kurokawa K, Nishimoto I. A
475 rescue factor abolishing neuronal cell death by a wide spectrum of familial Alzheimer's
476 disease genes and Abeta. *Proc Natl Acad Sci USA* 98: 6336-6341, 2001a. Erratum in:
477 *Proc Natl Acad Sci USA* 98, 2001. doi: 10.1073/pnas.101133498.
- 478 [7] Hashimoto Y, Ito Y, Niikura T, Shao Z, Hata M, Oyama F, Nishimoto I. Mechanisms of
479 neuroprotection by a novel rescue factor humanin from Swedish mutant amyloid
480 precursor protein. *Biochem Biophys Res Commun* 283: 460-468, 2001b.
481 doi:10.1006/bbrc.2001.4765.
- 482 [8] Lee C, Zeng J, Drew BG, Sallam T, Martin-Montalvo A, Wan J, Kim SJ, Mehta H,
483 Hevener AL, de Cabo R, Cohen P. The mitochondrial-derived peptide MOTS-c
484 promotes metabolic homeostasis and reduces obesity and insulin resistance. *Cell*
485 *Metab* 21: 433-454, 2015. doi:10.1016/j.cmet.2015.02.009.
- 486 [9] Cobb LJ, Lee C, Xiao J, Yen K, Wong RG, Nakamura HK, Mehta HH, Gao Q, Ashur C,
487 Huffman DM, Wan J, Muzumdar R, Barzilai N, Cohen P. Naturally occurring

488 mitochondrial-derived peptides are age-dependent regulators of apoptosis, insulin
489 sensitivity, and inflammatory markers. *Aging* 8: 786-809, 2016.
490 doi:10.18632/aging.100943.

491 [10] Lee C, Yen K, Cohen P. Humanin: a harbinger of mitochondrial-derived peptides?
492 *Trends Endocrinol Metab* 24: 222-228, 2013. doi:10.1016/j.tem.2013.01.005.

493 [11] Kim SJ, Xiao J, Wan J, Cohen P, Yen K. Mitochondrially derived peptides as novel
494 regulators of metabolism. *J Physiol* 595: 6613-6621, 2017. doi:10.1113/JP274472.

495 [12] Capt C, Passamonti M, Breton S. The human mitochondrial genome may code for
496 more than 13 proteins. *Mitochondrial DNA A DNA Mapp Seq Anal* 27: 3098-3101,
497 2016. doi:10.3109/19401736.2014.1003924.

498 [13] Logan IS. Pseudogenization of the Humanin gene is common in the mitochondrial DNA
499 of many vertebrates. *Zoological Res* 38:198-202, 2017. doi: 10.24272/j.issn.2095-
500 8137.2017.049

501 [14] Ladoukakis E, Pereira V, Magny EG, Eyre-Walke, A, Couso JP. Hundreds of putatively
502 functional small open reading frames in *Drosophila*. *Genome Biol* 12: R118, 2011.
503 doi:10.1186/gb-2011-12-11-r118.

504 [15] Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high
505 throughput. *Nucleic Acids Res* 32: 1792-1797, 2004. doi :10.1093/nar/gkh340.

506 [16] Kumar S, Stecher G, Tamura K. MEGA7: Molecular Evolutionary Genetics Analysis
507 Version 7.0 for Bigger Datasets. *Mol Biol Evol* 33: 1870-1874, 2016.
508 doi:10.1093/molbev/msw054.

509 [17] Vingron M, Brazma A, Coulson R, Van Helden J, Manke T, Palin K, Sand O, Ukkonen
510 E. Integrating sequence, evolution and functional genomics in regulatory genomics.
511 *Genome Biology* 10: 202-209, 2009. doi:10.1186/gb-2009-10-1-202.

512 [18] Rozas J. DNA sequence polymorphism analysis using DnaSP. *Methods Mol Biol* 537:
513 337-350, 2009. doi:10.1007/978-1-59745-251-9_17.

514 [19] Jukes TH, Cantor CR. *Evolution of protein molecules*. In: *Mammalian Protein*
515 *Metabolism*, edited by Munro HN, Academic Press, New York, 1969, p.21-132.
516 doi:10.1016/B978-1-4832-3211-9.50009-7.

517 [20] Lynch M, Crease TJ. The analysis of population survey data on DNA sequence
518 variation. *Mol Biol Evol* 7: 377-394, 1990. doi:10.1093/oxfordjournals.molbev.a040607.

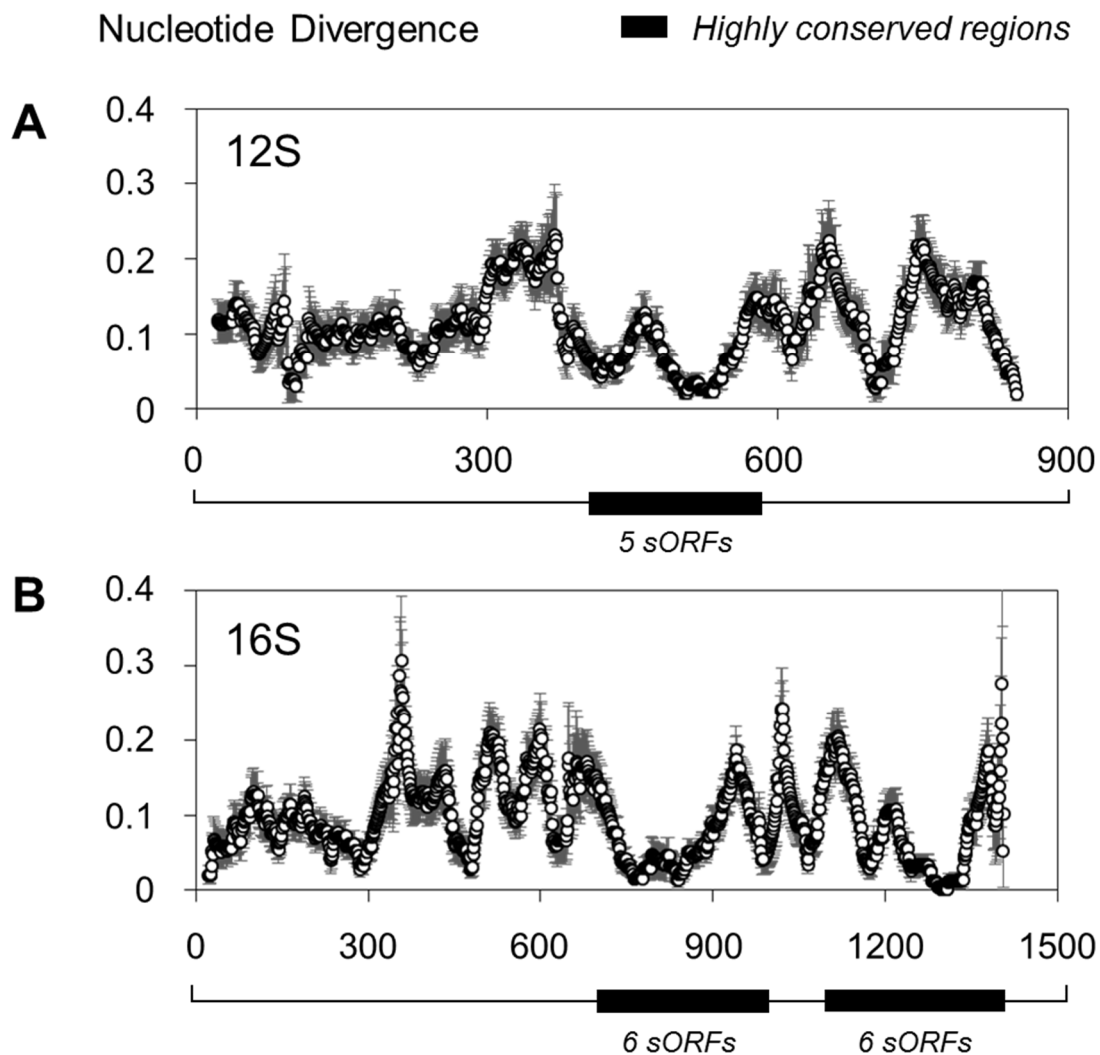
519 [21] Nei M. *Molecular Evolutionary Genetics*. New York: Columbia University Press, 1987.
520 doi:10.1002/ajpa.1330750317.

521 [22] Nei M, Gojobori T. Simple methods for estimating the numbers of synonymous and
522 nonsynonymous nucleotide substitutions. *Mol Biol Evol* 3: 418-426, 1986.
523 doi:10.1093/oxfordjournals.molbev.a040410.

- 524 [23] Nei M, Miller JC. A simple method for estimating average number of nucleotide
525 substitutions within and between populations from restriction data. *Genetics* 125: 873-
526 879, 1990.
- 527 [24] Osawa S, Jukes TH, Watanabe K, Muto A. Recent evidence for evolution of the genetic
528 code. *Microbiol Rev* 56: 229-264, 1992. doi:0146-0749/92/010229-36\$02.00/0.
- 529 [25] Watterson GA. On the number of segregating sites in genetical models without
530 recombination. *Theor Pop Biol* 7: 256-276, 1975. doi:10.1016/0040-5809(75)90020-9.
- 531 [26] Roussel D, Boël M, Mortz M, Romestaing C, Duchamp C, Voituron Y. Threshold effect
532 in the H₂O₂ production of skeletal muscle mitochondria during fasting and refeeding. *J*
533 *Exp Biol* 222, jeb196188, 2019. doi:10.1242/jeb.196188.
- 534 [27] Baker GC, Smith JJ, Cowan DA. Review and re-analysis of domain-specific 16S
535 primers. *J Microbiol Methods* 55: 541-555, 2003. doi:10.1016/j.mimet.2003.08.009.
- 536 [28] Wang Y, Tian RM, Gao ZM, Bougouffa S, Qian PY. Optimal eukaryotic 18S and
537 universal 16S/18S ribosomal RNA primers and their application in a study of symbiosis.
538 *PLoS One* 9, 2014. doi:10.1371/journal.pone.0090053.
- 539 [29] Gutell RR, Larsen N, & Woese CR. Lessons from an evolving rRNA: 16S and 23S
540 rRNA structures from a comparative perspective. *Microbiological Reviews*, 58(1), 10-
541 26, 1994.
- 542 [30] Fickett JW. ORFs and genes: how strong a connection? *J Comput Biol* 2: 117-123,
543 1995. doi:10.1089/cmb.1995.2.117.
- 544 [31] Guigó R, Fickett JW. Distinctive sequence features in protein coding genic non-coding,
545 and intergenic human DNA. *J Mol Biol* 253: 51-60, 1995. doi:10.1006/jmbi.1995.0535.
- 546 [32] Caricasole A, Bruno V, Cappuccio I, Melchiorri D, Copani A, Nicoletti F. A novel rat
547 gene encoding a Humanin-like peptide endowed with broad neuroprotective activity.
548 *FASEB J* 16: 1331-1333, 2002. doi:10.1096/fj.02-0018fje.
- 549 [33] Kariya S, Hirano M, Furiya Y, Ueno S. Effect of humanin on decreased ATP levels of
550 human lymphocytes harboring A3243G mutant mitochondrial DNA. *Neuropeptides* 39:
551 97-101, 2005. doi:10.1016/j.npep.2004.11.004.
- 552 [34] Thummasorn S, Shinlapawittayatorn K, Khamseekeaw J, Jaiwongkam T, Chattipakorn
553 SC, Chattipakorn N. Humanin directly protects cardiac mitochondria against
554 dysfunction initiated by oxidative stress by decreasing complex I activity. *Mitochondrion*
555 38: 31-40, 2018. doi:10.1016/j.mito.2017.08.001.
- 556 [35] Kozak M. An analysis of 5'-noncoding sequences from 699 vertebrate messenger
557 RNAs. *Nucleic Acids Research* 15, 1987. doi:10.1093/nar/15.20.8125.
- 558 [36] Kozak, M. Pushing the limits of the scanning mechanism for initiation of translation.
559 *Gene* 299(1-2), 1-34, 2002. doi:10.1016/S0378-1119(02)01056-9.

- 560 [37] Noderer WL, Flockhart RJ, Bhaduri A, Diaz de Arce AJ, Zhang J, Khavari PA, & Wang
561 CL. Quantitative analysis of mammalian translation initiation sites by FACS-seq.
562 *Molecular Systems Biology* 10(8), 748, 2014. doi:10.15252/msb.20145136.
- 563 [38] Park J, Kim J, Park W. Environmental stress to Ethiopia indigenous chicken breeds
564 induce immune signaling response among transcriptomic changes in heart, breast
565 muscle and spleen tissues. Public on Sep. 13, 2018. Available in GEO DataSets
566 (NCBI), ID of accession: GSE119387.
- 567 [39] Wang J. Transcriptomic Analysis between Normal and High-intake Feeding Geese
568 Provides Insight into Adipose Deposition and Susceptibility to Fatty Liver in Migratory
569 Birds. Public on Jan. 01, 2019. Available in GEO DataSets (NCBI), ID of accession:
570 GSE119421.
- 571 [40] Nei M, Kumar S. *Molecular Evolution and Phylogenetics*, Oxford University Press,
572 2000. doi:10.1046/j.1365-2540.2001.0923a.x.
- 573 [41] Krejcova G, Patocka J, & Slaninova J. Effect of humanin analogues on experimentally
574 induced impairment of spatial memory in rats. *Journal of Peptide Science*, 10(10), 636–
575 639, 2004. doi:org/10.1002/psc.569.
- 576 [42] Hashimoto Y, Niikura T, Ito Y, Sudo H, Hata M, Arakawa E, Nishimoto I. Detailed
577 characterization of neuroprotection by a rescue factor humanin against various
578 Alzheimer's disease-relevant insults. *The Journal of Neuroscience : The Official*
579 *Journal of the Society for Neuroscience* 21(23), 9235–45, 2001c.
580 doi:10.1523/JNEUROSCI.21-23-09235.2001.
- 581 [43] Wang D, Zhang S, He F, Zhu J, Hu S, Yu J. How do variable substitution rates
582 influence Ka and Ks calculations? *Genomics Proteomics Bioinformatics* 7: 116-27,
583 2009. doi:10.1016/S1672-0229(08)60040-6.
- 584 [44] Wang D, Liu F, Wang L, Huang S, Yu J. Nonsynonymous substitution rate (Ka) is a
585 relatively consistent parameter for defining fast-evolving and slow-evolving protein-
586 coding genes. *Biology Direct* 6-13, 2011. doi:10.1186/1745-6150-6-13.
- 587 [45] Xing Y, Lee C. Can RNA selection pressure distort the measurement of Ka/Ks? *Gene*
588 370: 1-5, 2006. doi:10.1016/j.gene.2005.12.015.
- 589 [46] Alekhina OM, & Vassilenko KS. Translation initiation in eukaryotes: Versatility of the
590 scanning model. *Biochemistry (Moscow)*, 77(13), 1465–1477, 2012.
591 <http://doi.org/10.1134/S0006297912130056>.
- 592 [47] Kochetov AV. Alternative translation start sites and hidden coding potential of
593 eukaryotic mRNAs. *BioEssays*, 30(7), 683–691. doi:10.1002/bies.20771.
- 594
- 595

596
597
598



599

600

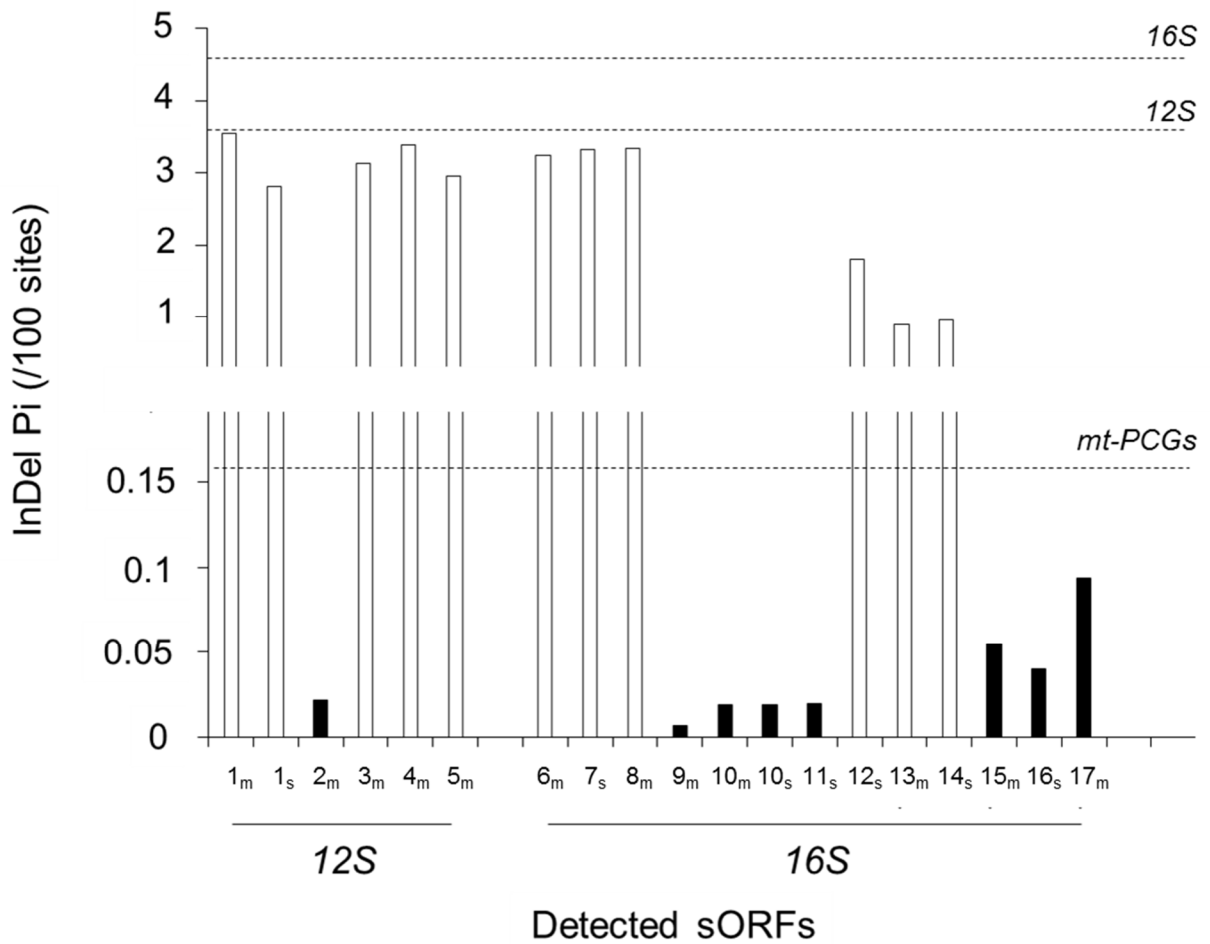
601 **Fig. 1.** Nucleotide divergences along 12S (**A**) and 16s (**B**) avian rDNAs (552 avian species).
602 Divergences were calculated for all positions without InDel events from the alignment of 12S
603 or 16S rDNAs. Values are the mean divergences \pm SEM of a 45 nucleotides sliding window.
604 Highly conserved regions identified in 12S or 16S rDNA (p -value $< 10^{-6}$) are represented by
605 dark rectangles. Note that other regions were significantly conserved ($p < 0.05$) (12S:
606 positions 213-300; 16S: positions 51-155 and 223-317). The consensus sequences of highly
607 conserved regions were used to detect highly conserved sORFs (≥ 45 nucleotides) in 12S
608 ($n=5$) and 16S ($n=11$) rDNAs by referring to mitochondrial and/or standard genetic codes.

609

610

611

612



613

614

615 **Fig. 2.** InDel polymorphism rates (InDel Pi) of detected sORFs in 12S or 16S rDNAs referring to either mitochondrial (m) or standard (s) genetic code. InDel Pi values calculated for the entire 12S or 16S rDNAs and for the concatenate of the thirteen mt-PCGs are indicated by 616 617 618 619 620 621 dotted lines. White bars represent sORFs with InDel Pi values close to those measured for the entire 12S (3.68%) or 16S (4.6%) rDNAs. Dark bars represent sORFs with InDel Pi values lower than that calculated for the concatenate of the 13 mt-PCGs (0.16%). Note that the y axis was broken and two scales are used for clarity.

622

623

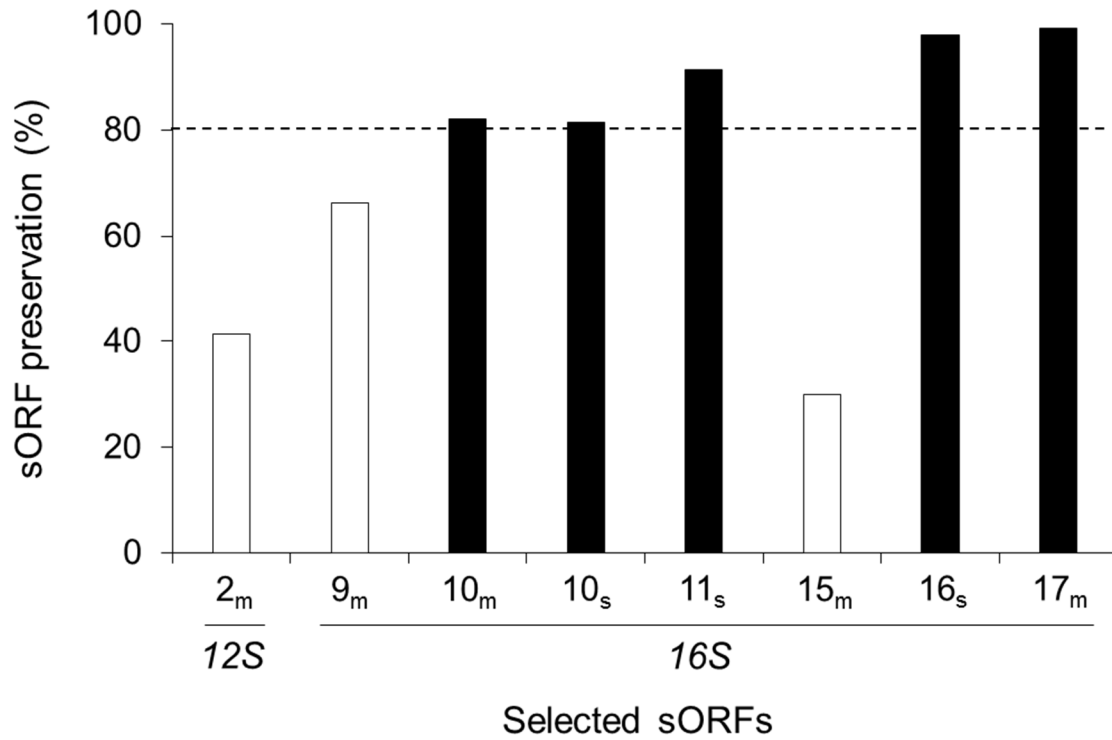
624

625

626

627

628



629

630

631 **Fig. 3.** sORF preservation of seven selected sORFs in bird species. The percentage of
632 preservation corresponds to the proportion of bird species presenting start and stop codons
633 in frame with initially detected consensus sORFs and flanking a region of ≥ 45 nucleotides.
634 Dark histograms represent the four sORFs preserved in more than 80% of bird species. Note
635 that sORF 9_{ms} is detected by using both mitochondrial and standard genetic codes.

636

637

638

639

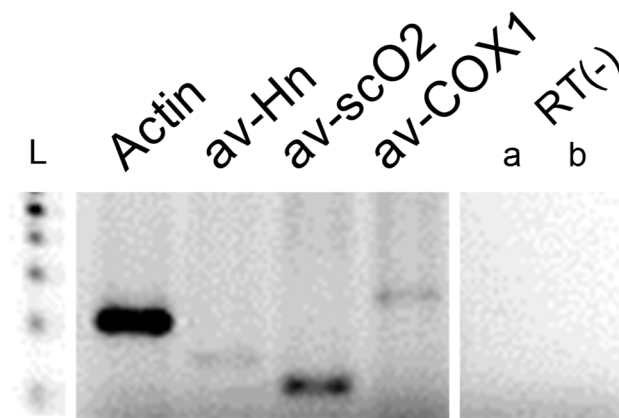
640

641

642

643

644

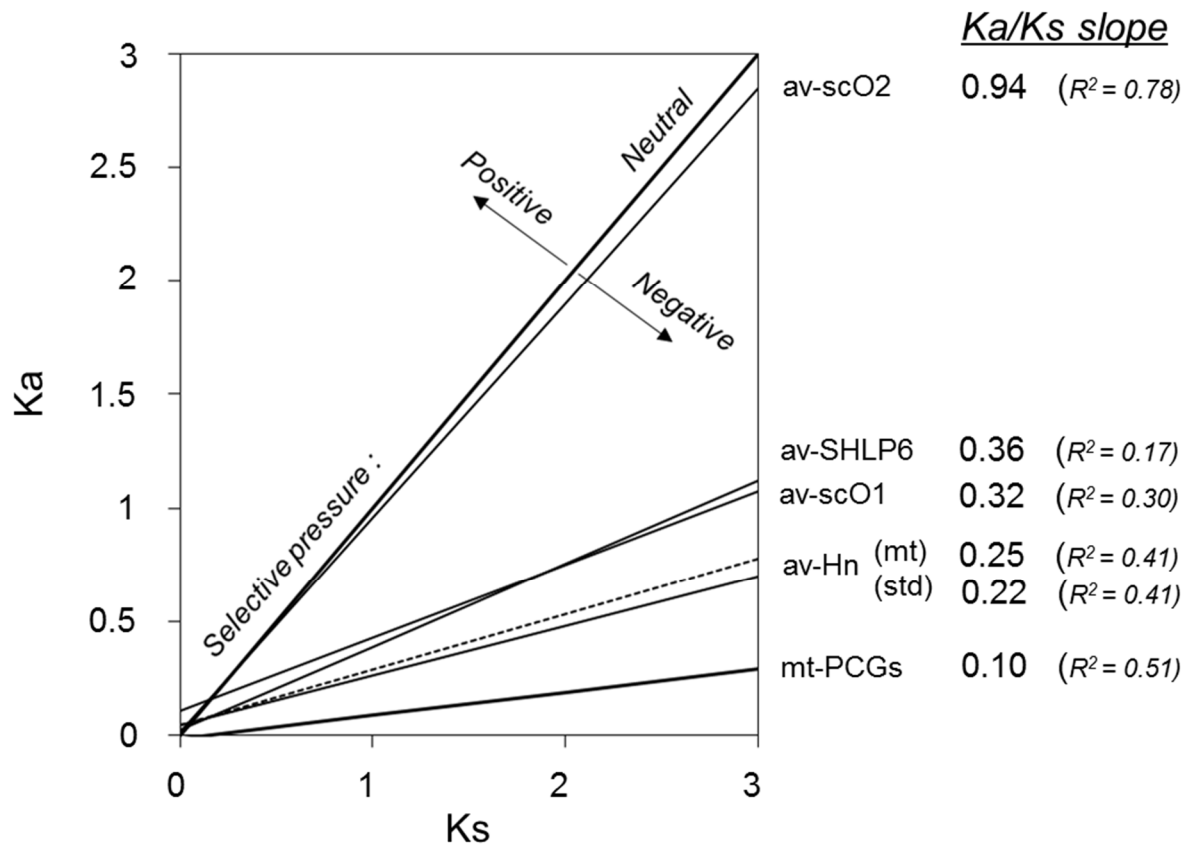


645

646

647 **Fig. 4.** RT-PCR analysis of the expression of av-Hn and av-scO2 in pectoralis muscle from a
648 6-wk-old muscovy duckling. RT-PCR reactions were performed using mRNAs that were
649 purified from total RNA extracted from muscle tissue with Dynabeads® mRNA purification kit
650 after DNase I digestion to remove traces of DNA and 30 cycles of PCR amplification were
651 used. Expected sizes were 155 bp for av-Hn and 127 bp for av-scO2. As positive controls,
652 avian actin (225 bp), encoded by a nuclear gene, and avian cytochrome c oxidase subunit 1
653 (av-Cox1, 230 bp), encoded by a mitochondrial gene, were used. As negative controls, PCR
654 amplifications of RNA not submitted to reverse transcription (RT(-)) were used with primers
655 for av-Hn (a) or av-scO2 (b) and poly(dT) to verify that only polyadenylated transcripts were
656 detected with RT-PCR. A 100 bp ladder (L) is shown.

657



658

659

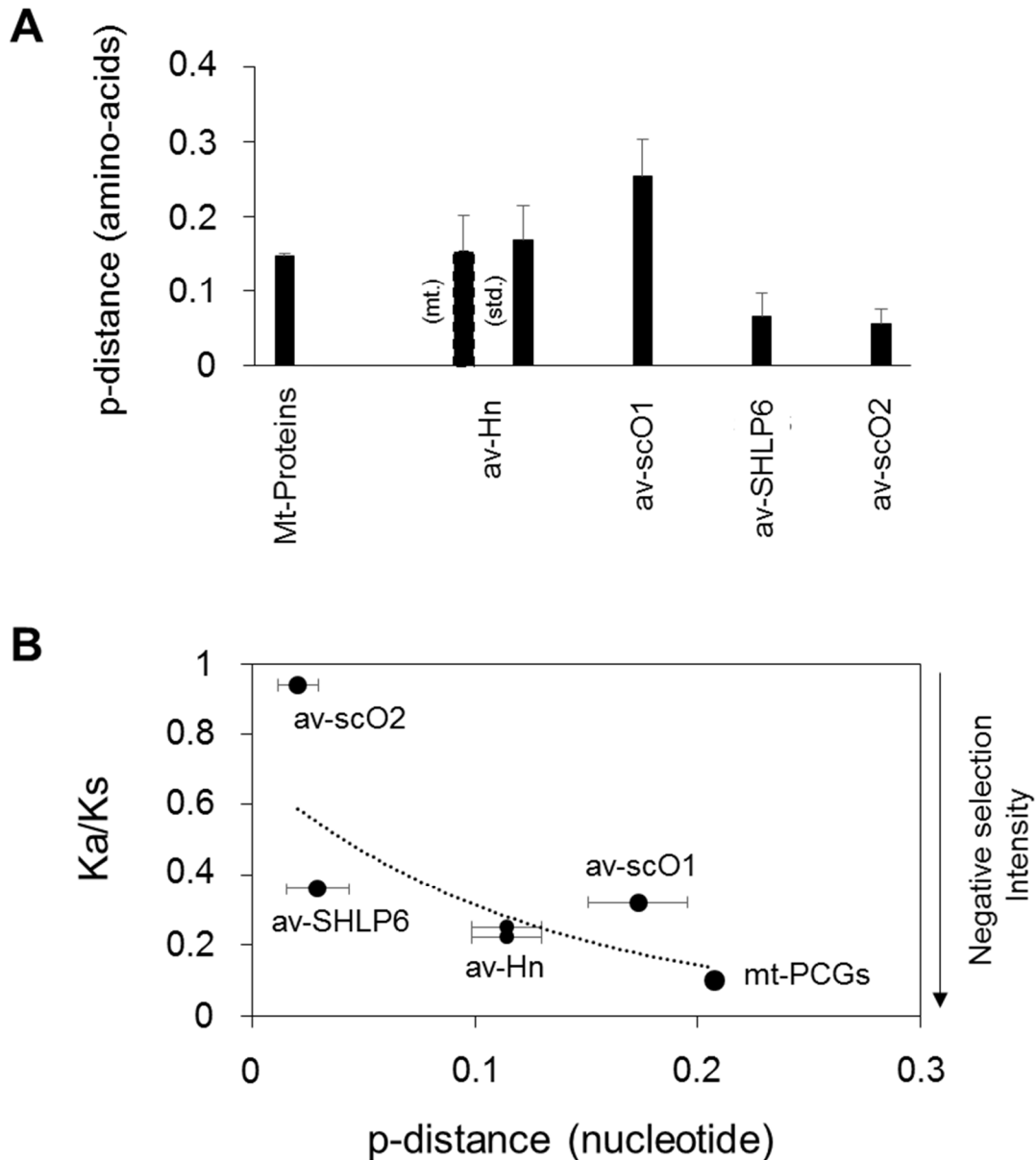
660 **Fig. 5.** Relation between Ka and Ks for the selected avian sORFs and avian mitochondrial
 661 protein coding genes (mt-PCGs). Trend lines were drawn for each sORF and for the
 662 concatenate of avian mt-PCGs and their respective slope (Ka/Ks) and R squared (R^2) values
 663 are indicated. These plots were obtained for each DNA region by calculating Ka and Ks
 664 values for pairwise comparisons among all studied bird species (552 bird species, 152 352
 665 pairwise comparisons per plot). The slope of trend lines is a good indicator of the selective
 666 pressure imposed on the nucleotide region, by indicating either positive ($Ka > Ks$), neutral
 667 ($Ka = Ks$) or negative selection ($Ka < Ks$).

668

669

670

671



672

673 **Fig. 6.** Conservation of a-a sequence for the proteins / peptides of interest (**A**) and relation
 674 between the rate of substitution and the negative selection intensity in birds. Proteins /
 675 peptides of interest were those potentially encoded by the selected sORFs and the
 676 concatenate of the known proteins encoded by mitochondrial protein-coding genes (mt-
 677 PCGs). The p-distance among all 552 studied bird species was calculated by using either
 678 translated peptide (**A**) or nucleotide (**B**) sequences. Variances were estimated through 100
 679 bootstrap replications of each calculation. The relation between the rate of substitution and
 680 the negative selection intensity (**B**) was drawn by plotting Ka/Ks ratio as a function of
 681 nucleotide p-distance value for each region. The smaller the Ka/Ks ratio, the higher the
 682 relative intensity of negative selection.

683

684 **Table 1**

685 Comparison of avian selected sORFs with potential human counterparts.

sORF	Position	Gen. Code	Strand	Species	a-a sequence	Identification
10_{ms}	886-960	Mt/Std	+	Birds	MAK RGLNCLLQVI <u>SEIDL</u> PVQKQG*	av-Hn (<i>avian orthologue of Humanin</i>)
				Human	MAP RGFSCLLLLT <u>SEIDL</u> PVKRRA*	Humanin (<i>Hashimoto et al, 2001</i>)
16_s	1238-1300	Std	+	Birds	MLD QDILMVQP LLRVRLFND *	av-SHLP6 (<i>avian orthologue of SHLP6</i>)
				Human	MLD QDIPMVQP LLKVRLFND *	SHLP6 (<i>Cobb et al, 2016</i>)
11_s	913-966	Std	-	Birds	MCVFLLLYREIN FTDYL *	av-scO1 (<i>avian strongly conserved ORF1</i>)
				Human	AVLC PPL HGQV NFTG *K*	No homologous sORF
17_m	1274-1351	Mt	-	Birds	ME <u>TDLDC</u> SGLNSDHVGLLIVEQTNP*	av-scO2 (<i>avian strongly conserved ORF2</i>)
				Human	ME <u>TDL</u> DYSGLNSDHVGL*	Homologous sORF not yet identified

686

687 Bird consensus sequences and their human (NC_012920.1, Pos.1671:3229) homologous regions were aligned by codons and translated into a-
688 a sequences using either mitochondrial and/or standard genetic (gen.) codes. Putative orthologs of human peptides were renamed accordingly.
689 sORFs with no human counterpart so far were also renamed. The position of each sORF in consensus sequence and putative encoding strand
690 are specified. Bold letters indicate conserved a-a in birds and human. Analysis of putative peptides with the PROSITE database
691 (<https://prosite.expasy.org/>) led to identify functional domains in av-Hn and av-scO2 sequences such as putative casein kinase II
692 phosphorylation sites (underlined bold italic letters) that are fully conserved between birds and human in either av-Hn or av-scO2.

693