1          **Article**

2

3      **Contrasted gene decay in subterranean vertebrates: insights from**

4      **cavefishes and fossorial mammals**

5

6      Maxime Policarpo[1], Julien Fumey[‡,1], Philippe Lafargeas[1], Delphine Naquin[2], Claude

7      Thermes[2], Magali Naville[3], Corentin Dechaud[3], Jean-Nicolas Volff[3], Cedric Cabau[4],

8      Christophe Klopp[5], Peter Rask Møller[6], Louis Bernatchez[7], Erik García-Machado[7,8], Sylvie

9      Rétaux[*,9] and Didier Casane[*,1,10]

10

11     [1] Université Paris-Saclay, CNRS, IRD, UMR Évolution, Génomes, Comportement et
12     Écologie, 91198, Gif-sur-Yvette, France.
13     [2] Institute for Integrative Biology of the Cell, UMR9198, FRC3115, CEA, CNRS, Université
14     Paris-Sud, 91198 Gif-sur-Yvette, France.
15     [3] Institut de Génomique Fonctionnelle de Lyon, Univ Lyon, CNRS UMR 5242, Ecole
16     Normale Supérieure de Lyon, Université Claude Bernard Lyon 1, Lyon, France.
17     [4] SIGENAE, GenPhySE, Université de Toulouse, INRAE, ENVT, F-31326, Castanet Tolosan,
18     France.
19     [5] INRAE, SIGENAE, MIAT UR875, F-31326, Castanet Tolosan, France.
20     [6] Natural History Museum of Denmark, University of Copenhagen, Universitetsparken 15,
21     DK-2100 Copenhagen Ø, Denmark.
22     [7] Department of Biology, Institut de Biologie Intégrative et des Systèmes, Université Laval,
23     1030 Avenue de la Médecine, Québec City, Québec G1V 0A6, Canada.
24     [8] Centro de Investigaciones Marinas, Universidad de La Habana, Calle 16, No. 114 entre 1ra e
25     3ra, Miramar, Playa, La Habana 11300, Cuba.
26     [9] Université Paris-Saclay, CNRS, Institut des Neurosciences Paris-Saclay, 91190, Gif-sur-
27     Yvette, France.
28     [10] Université de Paris, UFR Sciences du Vivant, F-75013 Paris, France.
29
30
31     [‡] Present address: Human Genetics and Cognitive Functions, Institut Pasteur, CNRS UMR
32     3571, Université de Paris, Paris 15
33
34     * Corresponding authors: E-mails: sylvie.retaux@inaf.cnrs-gif.fr; didier.casane@egce.cnrs-
35     gif.fr.
36

## Abstract (241 words; max = 250)

Evolution sometimes proceeds by loss, especially when structures and genes become dispensable after an environmental shift relaxing functional constraints. Gene decay can serve as a read-out of this evolutionary process. Animals living in the dark are outstanding models, in particular cavefishes as hundreds of species evolved independently during very different periods of time in absence of light. Here, we sought to understand some general principals on the extent and tempo of decay of several gene sets in cavefishes. The analysis of the genomes of two Cuban species belonging to the genus *Lucifuga* provides evidence for the most massive loss of eye genes reported so far in cavefishes. Comparisons with a recently-evolved cave population of *Astyanax mexicanus* and three species belonging to the tetraploid Chinese genus *Sinocyclocheilus* revealed the combined effects of the level of eye regression, time and genome ploidy on the number of eye pseudogenes. In sharp contrast, most circadian clock and pigmentation genes appeared under strong selection. In cavefishes for which complete genomes are available, the limited extent of eye gene decay and the very small number of loss of function (LoF) mutations per pseudogene suggest that eye degeneration is never very ancient, ranging from early to late Pleistocene. This is in sharp contrast with the identification of several eye pseudogenes carrying many LoF mutations in ancient fossorial mammals. Our analyses support the hypothesis that blind fishes cannot thrive more than a few millions of years in cave ecosystems.

*Key words:* cavefishes, eye genes, pseudogenization, machine learning, relaxed selection, molecular dating.

## Introduction (791 words)

The evolution of organisms confronted to drastic environmental shifts results in sometimes profound phenotypic changes. Constructive evolution involved in adaptation to new environments, and relying on novelties at phenotypic and genetic levels, has drawn much interest. Nevertheless, it becomes evident that regressive evolution, which is often non adaptive and which occurs by loss of structures and functions and corresponding genes, accounts for a non-negligible part of the evolutionary process (Lahti, et al. 2009; Albalat and Cañestro 2016). Here, we sought to better understand the modalities, extent, tempo and limits of molecular decay of several light-related genetic systems in subterranean vertebrates. It has been shown that several independent lineages of obligate fossorial mammals with degenerated eyes have lost many genes involved in visual perception (Kim, et al. 2011; Emerling and Springer 2014; Fang, Nevo, et al. 2014; Fang, Seim, et al. 2014; Emerling 2018). Cave vertebrates which are essentially cavefishes are other outstanding models to tackle these issues (Culver and Pipan 2009). However, the molecular decay of genes has not been surveyed at a genome-wide scale in relevant cavefish species. On the one hand, in the reference genome of *A. mexicanus* cavefish, no or only a couple of pseudogenes have been found among sets of genes which are eye specific, involved in the circadian clock, or else related to pigmentation (Protas, et al. 2006; Beale, et al. 2013; McGaugh, et al. 2014). Such maintenance of a very high proportion of functional genes most likely results from a very recent origin, no earlier than in the late Pleistocene, of cave populations (Fumey, et al. 2018). On the other hand, in the genomes of three fishes belonging to the genus *Sinocyclocheilus* (*S. grahami* which is a surface fish with large eyes, *S. anshuiensis* which is a blind cavefish and *S. rhinocerous* which is a small-eyed cavefish) many LoF mutations were found (Yang, et al. 2016), but their tetraploid genomes hampered the identification of those mutations that fixed

87    in relation to the surface to cave shift. Indeed, after a whole-genome duplication (WGD), the

88    pair of paralogs resulting from this process (ohnologs) are most often redundant and one of

89    them can be pseudogenized without reducing fitness. Accordingly, *S. grahami* carries eye

90    peudogenes like the eyeless *S. anshuiensis* and the small-eyed *S. rhinocerous,* but no thorough

91    analysis of differential gene losses in relation to the level of eye degeneration has been

92    performed (Yang, et al. 2016).

93    In order to examine the long term effect of life in caves on the molecular decay of large sets

94    of genes involved in various light-dependent biological processes, genomes of fishes evolving

95    in caves for a very long time and which did not undergo a recent WGD are required. Two

96    clades of cavefishes (cave brotulas from Bahamas and Cuba) were previously identified in the

97    genus *Lucifuga*, one comprising only blind cavefish species and the other only small-eyed

98    cavefish species (García-Machado, et al. 2011). As no close surface relative has been

99    identified up to now and large genetic distances were found between some species, within and

100   between clades, this genus of cavefishes is likely relatively ancient, and the last common

101   ancestor of extant species was probably a cave-adapted fish. We sequenced the genomes of

102   two Cuban cave brotulas: one specimen, belonging to *L. dentata*, was blind and depigmented,

103   the other one, belonging to *L. gibarensis*, had small eyes and was pigmented. The latter

104   species is a new species first identified as *Lucifuga* sp. 4 (García-Machado, et al. 2011) that

105   will be formally named and described in a forthcoming publication.

106   We searched for likely LoF mutations (*i.e.* STOP codon gains, losses of START and STOP

107   codons, losses of intron splice sites and small indels leading to frameshifts) and for several

108   signatures of relaxed selection on nonsynonymous mutations in genes: 1) uniquely expressed

109   in the eyes or coding for non-visual opsins, 2) involved in the circadian clock, 3) involved in

110   pigmentation. The contrasted patterns of pseudogenization found for the three categories of

111   genes indicate that eye genes are much less constrained than circadian clock and pigmentation

112    genes in caves. In *A. mexicanus* cavefish, despite only one eye gene carrying a LoF mutation

113    was found, using machine learning-based estimations of the deleterious impact of

114    nonsynonymous mutations implemented in MutPred2 (Pejaver, et al. 2017), we obtained

115    evidence that most if not all eye genes are under relaxed selection, but for a too short period

116    of time to allow the fixation of more than a few LoF mutations. In other cavefishes, more eye

117    pseudogenes were found and the level of gene decay depended on several factors such as the

118    time fishes have spent in the subterranean environment, their level of troglomorphy and the

119    level of ploidy of their genomes. Nevertheless, no eye genes with many LoF mutations were

120    found, in sharp contrast to highly degenerated eye genes identified in some fossorial

121    mammals, suggesting that eye degeneration in cavefishes is much more recent.

122

123    **Results**

124

125    Assembly of the draft genomes of two Cuban cave brotulas

126

127    First, the genome of a specimen of *Lucifuga dentata* (Bythitidae, Ophidiiformes) was

128    sequenced (see a photo in **supplementary fig. S1, Supplementary Material** online).

129    Assembly resulted in 52,944 scaffolds whose size sum up to 634 Mb, N50 = 119.6 kb (for

130    scaffold size distribution, see **supplementary fig. S2, Supplementary Material** online). This

131    genome size is consistent with those of three other genomes available (Malmstrøm, et al.

132    2017) and estimations of the genome size of other Ophidiiformes (Gregory 2019). To assess

133    the quality of the assembly, raw sequences were realigned to the assembly: 95% of the reads

134    realigned correctly resulting in a mean coverage of 134x. Then, the completeness of the

135    assembly was assessed using BUSCO with the Actinopterygii gene database (Kriventseva, et

136    al. 2015). Among 4,584 genes, 4,249 (92.7%) were found complete, 194 (4.2%) were

137  incomplete and 141 (3.1%) were missing. Using BUSCO with three other Ophidiiformes

138  genomes currently available (*Brotula barbata*, *Carapus acus* and *Lamprogrammus exutus*),

139  the genome of *Lucifuga dentata* appeared as the most complete (see **supplementary fig. S3,**

140  **Supplementary Material** online). Then, the genome of a specimen belonging to the small-

141  eyed *Lucifuga gibarensis* was sequenced (see a photo in **supplementary fig. S1,**

142  **Supplementary Material** online). As nuclear DNA sequence divergence is about 1%

143  between the two *Lucifuga* species, short reads of *L. gibarensis* were mapped on *L. dentata*

144  genome. The mean coverage was 84x, with 86% of the reads mapping on the genome.

145  Heterozygosity was estimated on raw Illumina reads using GenomeScope (Vurture, et al.

146  2017). The heterozygosity of *L. dentata* (0.1%) was lower than that of *L. gibarensis* (0.26%).

147

148  Assembly of a transcriptome of *L. dentata* and genome annotation

149

150  Based on mRNA extracted from the gonads, gills, heart and brain of *L. dentata*, a *de novo*

151  transcriptome assembly was obtained using Trinity (Grabherr, et al. 2011). Quality and

152  completeness assessment of this transcriptome were performed following Trinity guide.

153  Among 4,584 genes corresponding to the Actinopterygii gene database of BUSCO, 82.8%

154  were found complete (**supplementary fig. S3, Supplementary Material** online) and 92.31 %

155  of the reads were mapped back to the assembly with 84 % as proper pairs, which indicate an

156  overall good quality transcriptome. More on quality check can be found in **supplementary**

157  **fig. S4, Supplementary Material** online.

158  A combination of *de novo* predictions, RNA-seq evidence and protein alignments was used to

159  annotate the genome of *L. dentata* (see workflow in **supplementary fig. S5, Supplementary**

160  **Material** online). This resulted in 30,001 gene models with an average gene length of 9,693

161  bp and an average protein length of 435 amino acids. Among predicted genes, 23,524 had a

162    functional annotation with BLAST to the SwissProt/UniProt database and 21,558 genes were

163    detected with a functional domain by Interproscan. Annotation completeness was assessed

164    using BUSCO in protein mode; among 4,584 corresponding to the Actinopterygii gene

165    database of BUSCO, 87.4% were found complete, 6.5% incomplete and 6% missing

166    (**supplementary fig. S3, Supplementary Material** online). A homemade pipeline was used

167    to describe the repeat landscape of the genome of *L. dentata*. We found 16.3% of repeated

168    elements, among which 2.4% of LINEs and 0.4% of SINEs (**supplementary fig. S6,**

169    **Supplementary Material** online).

170

171    Delimitation and retrieving of eye, circadian clock and pigmentation genes

172

173    In zebrafish, *Danio rerio,* we identified 95 genes expressed only in the eyes or coding non-

174    visual opsins expressed in other organs (**fig. 1A, supplementary fig. S7** and **Data Supp 1,**

175    **Supplementary Material** online, and see Methods). In addition, we retrieved a list of 42

176    circadian clock genes (Li, et al. 2013) and 257 genes involved in pigmentation (Lorin, et al.

177    2018) (**fig. 1B** and **fig. 1C, supplementary Data_Supp1, Supplementary Material** online).

178    Using the program exonerate, homologs were retrieved from other fish genomes, that is five

179    cavefishes (*A. mexicanus* from Pachón cave, *L. dentata* and *L. gibarensis*, *S. anshuiensis* and

180    *S. rhinocerous*), close surface relatives (*A. mexicanus* and *Pygocentrus nattereri*, *Brotula*

181    *barbata*, *Carapus acus* and *Lamprogrammus exutus, S. grahami* and *C. carpio)* and a

182    distantly-related outgroup (*Lepisosteus oculatus*). Their phylogenetic relationships are shown

183    in **fig. 2**. Noteworthy, some genes have been duplicated in the terminal lineage leading to

184    zebrafish (used as a reference to establish the gene lists) and thus only one copy was expected

185    to be found in other fishes. On the other hand, gene duplications, gene deletions as well as

186    WGDs occurred in other lineages. Therefore, the number of genes retrieved is highly variable

187    among genomes (**fig. 3**).

188

189    ## Identification of LoF mutations

190

191    Genes sequences were classified as functional if found complete with no LoF mutation, as

192    pseudogene if complete and carrying at least one LoF mutation, and as truncated if incomplete

193    (the sequences can be found in **supplementary Data_Supp2, Supplementary Material**

194    online). Only the following LoF mutations were analyzed: gain of an internal STOP codon,

195    loss of the initiation codon, loss of the STOP codon, indel leading to a frameshift, mutations

196    at intron donor and acceptor sites. In the present study, incomplete genes were discarded as it

197    was difficult to know if they corresponded to sequencing gaps, assembly artefacts or true

198    large deletions. Using PCR to amplify missing exons, we estimated that 85% of the large

199    deletions in the *A. mexicanus* cavefish genome are artefacts (data not shown) - although some

200    large deletions such as in the gene *Oca2* (a pigmentation gene) are real (Protas, et al. 2006).

201    Other mutations in non-coding and coding sequences that could lead to a non-functional gene

202    were not searched for as they cannot be readily identified. For example, several in-frame indel

203    mutations were found in *A. mexicanus* but their functional consequences remained elusive

204    (Berning, et al. 2019). The numbers of pseudogenes reported hereafter are thus underestimates

205    of the true numbers of non-functional genes, but they nevertheless allowed comparative

206    analyses.

207

208    Eye pseudogenes: among the list of 95 zebrafish eye genes, 76 genes were retrieved from

209    *Lucifuga* genomes, 75 from *B. barbata*, 72 from *C. acus* and 73 from *L. exutus* (**fig. 3**,

210    **Supplementary fig. S7** and **Data_Supp1, Supplementary Material** online). Interestingly,

8

211    all these ophidiiforms seem to have lost long-wave sensitive (LWS) opsins. This loss is most

212    likely due to a gene deletion in their common ancestor living in deep ocean (**supplementary**

213    **fig. S8, Supplementary Material** online), in accordance with a report on the reduction of the

214    number of LWS genes in fishes living below 50 m (Lin, et al. 2017). While no eye

215    pseudogene was found in *B. barbata* or *C. acus* and only one in *L. exutus* (*gcap1*), 5

216    pseudogenes were identified in *L. gibarensis* and 19 pseudogenes in *L. dentata*. The non-

217    visual opsin *rgr1* was pseudogenized in the common ancestor of the two *Lucifuga* species, as

218    the same mutation (at a splice site of intron 4) was found in both genomes (**fig. 3** and

219    **supplementary Data_Supp2, Supplementary Material** online). Examination of the read

220    coverage of LoF mutations indicated that the specimen of *L. gibarensis* sequenced was

221    heterozygous for LoF mutations found at two different sites in the *gcap2* gene

222    (**supplementary table S1, Supplementary Material** online). In the transcriptome of *L.*

223    *dentata*, transcripts corresponding to 9 pseudogenes were found (3 non-visual opsins, 3

224    crystallins and 3 genes involved in the phototransduction pathway), while no transcripts were

225    found for 10 other pseudogenes (**supplementary table S1, Supplementary Material** online).

226    In those transcripts, all the LoF mutations identified at the genome level were present.

227    In agreement with a recent WGD, two copies (ohnologs) of most eye genes were retrieved

228    from the genomes of *Sinocyclocheilus* species (**fig. 3**, **supplementary fig. S7,**

229    **Supplementary Material** online). In the large-eyed *S. grahami*, about 10% of retrieved eye

230    genes were pseudogenized (18 / 173 genes carried at least a LoF mutation), to be compared to

231    19% (32 / 169) in the small-eyed *S. rhinocerous* and 28% (48 / 171) in the eyeless cavefish *S.*

232    *anshuiensis*. Only one pair of ohnologs were concomitantly pseudogenized in the eyed *S.*

233    *grahami* and the small-eyed *S. rhinocerous,* while seven pairs of ohnologs were

234    concomitantly pseudogenized in the blind *S. anshuiensis* (**fig.1, fig.3, supplementary fig. S7,**

235    **Data_Supp1, Supplementary Material** online). A STOP codon and a frameshift in *sws1*

236    were shared by the three *Sinocyclocheilus* species and *Cyprinus carpio*. A new STOP codon

237    and a new frameshift in this gene were shared by *Sinocyclocheilus* species, as well as a

238    mutation at the donor site of the third intron of *gc3*; *S. anshuiensis* and *S. grahami* shared a

239    frameshift in *crygm5* and a frameshift and a new STOP codon in *grk7b* (**fig. 3**).

240    In *A. mexicanus*, 86 genes were retrieved from the surface fish genome while 85 were

241    retrieved from the genome of the Pachón cavefish. Only one pseudogene was found in the

242    Pachón cavefish genome, which is due to a deletion of 11 bp in *pde6b* (**fig. 1** and **fig. 3**). The

243    examination of the automatic annotation of the gene allowed the identification of an erroneous

244    1 bp intron (ENSAMXG00000000290, Ensembl 91), restoring the coding frame. Noteworthy,

245    we also confirmed using PCR that two large deletions occurred in the Pachón cavefish, one

246    removing *opn8b* and the last 3 exons of *opn8a,* the other eliminating 2 exons of *rgr2*.

247    However, these genes were not included in the list of pseudogenes, according to the restrictive

248    definition of an identifiable pseudogene used in the present study.

249    In summary, while no or very few eye genes are pseudogenized in surface fishes and *A.*

250    *mexicanus* cavefish, more eye pseudogenes were found in other cavefishes, up to 25% in *L.*

251    *dentata*.

252

253    Circadian clock pseudogenes: based on a literature survey, 42 genes involved in the circadian

254    clock in *Danio rerio* were identified (Li, et al. 2013) and retrieved from other fish genomes.

255    On the one hand, no pseudogene was found among 36 genes retrieved from *Lucifuga* genomes

256    and 38 genes identified from *Astyanax* genomes. On the other hand, 5, 15 and 9 pseudogenes

257    were identified among 80, 83 and 81 genes retrieved from the genomes of *S. grahami* (eyed),

258    *S. rhinocerous* (small-eyed) and *S. anshuiensis* (blind)*,* respectively. Both ohnologs of *cry-*

259    *dash* were independently pseudogenized in *S. rhinocerous* and *S. anshuiensis*, a gene also

260    pseudogenized in the Somalian cavefish *P. andruzzii*, Three other pair of ohnologs (*cry1b*,

261   *cry2a* and *per2*) carried LoF mutations in *S. rhinocerous*, *per2* being also non-functional in *P.*

262   *andruzzii*. These data suggest that the circadian clock has most likely been lost in *S.*

263   *rhinocerous* whereas this loss is less strongly supported in the case of *S. anshuiensis* (**fig. 1**

264   and **fig. 3**).

265

266   <u>Pigmentation genes</u>: based on a literature survey, 257 genes involved in pigmentation in

267   *Danio rerio* were identified (Lorin, et al. 2018) and retrieved from other fish genomes. Very

268   few pseudogenes were found among 237 pigmentation genes in *Lucifuga* genomes, that is 8

269   LoF mutations in *L. dentata* and 7 in *L. gibarensis*. While *smtla* and *myo7ab* seem to have

270   been lost independently in the two lineages, a STOP codon and an insertion is shared in

271   *adamts20*. The number of pseudogenes in these cavefishes does not seem to depart from those

272   found in some surface relatives, as 6 pseudogenes were identified among 230 pigmentation

273   genes in *Lamprogrammus exutus* (**supplementary Data_Supp1, Supplementary Material**

274   online). Among *Sinocyclocheilus* species, only 3% (15/484) of pseudogenes were found in *S.*

275   *grahami* while 6% (28/490) were found in *S. rhinocerous* and 7% (35/487) in *S. anshuiensis*

276   (**fig. 3**). Thus, after the WGD, the retention of pigmentation genes seems to have been much

277   higher than among eye genes in the two cavefishes but also in the surface fish (compare to

278   10%, 19% and 28% of eye pseudogenes, respectively). Such high percentage of retention of

279   pigmentation genes has been found also after the Salmonid-specific WGD (Lorin, et al. 2018).

280   Strikingly, while no pair of ohnologs was found pseudogenized in *S. grahami*, the same two

281   pairs of ohnologs (*gch2* and *pmelb*) were independently pseudogenized in *S. anshuiensis* and

282   *S. rhinocerous*. The very small number of pseudogenes and the independent pseudogenization

283   of the same genes in these two species suggest that only a limited subset of genes involved in

284   pigmentation can be lost in these cavefishes.

285    In *A. mexicanus*, 2 pseudogenes were found among 249 pigmentation genes: *mc1r* which has

286    already been reported in the literature (Gross, et al. 2009) and which is also pseudogenized in

287    the Chinese cavefish *Oreonectes daqikongensis* (Liu, et al. 2019), and *tyrp1a* that is mis-

288    annotated in ensembl (ENSAMXG00000021619, Ensembl 91).

289

290    Types of LOF mutations and their distribution along pseudogenes

291

292    A total of 118 mutations to a STOP codon, 148 frameshifts (84 deletions and 64 insertions), 5

293    STOP codon losses, 13 START codon losses and 40 intron splice site losses were identified in

294    our dataset (**fig. 4**) (see **supplementary Data_Supp1**, **Supplementary Material** online for a

295    detailed description of the number of LoF mutations in each gene set). Most frameshifts were

296    the results of very small deletions or insertions (1 or 2 bp) while a few of them were indels

297    involving a larger number of nucleotides (**fig. 4C, supplementary fig. S10**, **Supplementary**

298    **Material** online). The largest deletion was 83 bp long in *zic2b* (pigmentation gene) of *S.*

299    *rhinocerous* and the largest insertion was 20 bp long in *opn7b* (eye gene) of *S. anshuiensis*. In

300    order to test if LoF mutations were distributed randomly along the genes, that is they were not

301    clustered at the 3' end of the genes where their deleterious effect could be low, we computed

302    the effective segment size generated by LoF mutations and compared this value with

303    simulations of random distributions of mutations along genes. We found that premature STOP

304    codons and frameshifts are distributed randomly along coding sequences (for more details, see

305    Materials and Methods and **supplementary fig. S9**, **Supplementary Material** online).

306    We next tested if the relative frequencies of the different types of LoF mutations (*i.e.* STOP

307    codon gains, losses of START and STOP codons, losses of intron splice sites and small indels

308    leading to frameshifts) were those expected under the neutral model, that is if their relative

309    frequencies were proportional to their probabilities of occurrence. Taking into account a

310  nucleotide mutation rate (μ), observed transition/transversion ratio and codon frequencies, the

311  rate of mutation to a new STOP codon was $\mu_{stop} = 0.036\mu$. Based on the ratio of

312  frameshift/STOP mutations, we estimated the rate of indels leading to frameshifts ($\mu_{frameshift}$)

313  as 148/118 x $\mu_{stop} = 0.05\mu$. The rate of mutations in splice acceptor or donor sites was

314  estimated as: 4 x (number of introns) / Σ (CDS length) x μ, *i.e.* 4 x 34175 / 6154965 x μ =

315  0.022μ (where 34175 is the number of introns and 6154965 is the number of bases identified

316  in the 3625 genes retrieved from the genomes of two *Lucifuga* species, three *Sinocyclocheilus*

317  species and two genomes of *Astyanax mexicanus*). The rate of START codon loss was

318  estimated as: 3 x (number of genes) / Σ (CDS length) x μ, *i.e.* 3 x 3625 / 6154965 x μ =

319  0.0018μ. The rate of STOP codon loss was estimated as: 3 x (number of genes) / Σ (CDS

320  length) x 0.85 x μ, *i.e.* 3 x 3625 / 6154965 x 0.85 x μ = 0.0015μ (where 0.85 is the probability

321  that a mutation in a STOP codon leads to a sense codon). The observed distribution of LoF

322  mutations fitted well with those expected, either taking the three datasets together (**fig. 4B,**

323  **supplementary fig. S11** and **Data_Supp1**, **Supplementary Material** online) or each dataset

324  individually. It was similar to the distribution found in eye pseudogenes of subterranean

325  mammals (**supplementary fig. S11 and supplementary Data_Supp1**, **Supplementary**

326  **Material** online). These results suggested that the number LoF mutations of each type is

327  proportional to its probability of occurrence.

328

329  Estimation of the number of effectively neutral eye genes based on the

330  distribution of LoF mutations per pseudogene in cave brotulas

331

332  Among pseudogenes, some accumulated more than one LoF mutation, but in most of the

333  cases only one LoF mutation was found (**supplementary fig. S7 and Data_Supp1**,

334   **Supplementary Material** online). In order to test if the whole set, or only a subset, of eye

335   genes is free to accumulate LoF mutations, we compared the distribution of the number of

336   LoF mutations per pseudogene with those expected under these different hypotheses.

337   Expected distributions were obtained using either a simple analytical model assuming that all

338   genes have the same probability to fix a LoF mutation, or a more complex model that takes

339   into account that different genes do not have the same probability to fix a LoF mutation

340   because they have different length and they do not contain the same number of introns. In the

341   latter case, the computation of expected distribution was based on simulations. Very similar

342   expected distributions were obtained with both approaches. This analysis could be performed

343   only with *Lucifuga* species, as only one LoF mutation was found in *Astyanax mexicanus* and a

344   WGD allowed LoF mutations to reach fixation in a *Sinocyclocheilus* species with large

345   functional eyes such as *S. grahami*.

346   In *L. dentata*, 22 LoF mutations were distributed among 19 eye pseudogenes. More precisely,

347   among the 76 genes retrieved, there were 57 genes without LoF mutation, 16 with 1 mutation,

348   and 3 with 2 mutations (**fig. 3** and **supplementary fig. S7**, **Supplementary Material** online).

349   This distribution was compared with expected distributions obtained for different numbers of

350   neutral genes ranging from 19 to 76 (**fig. 5A**). The best fit between the observed and expected

351   distribution was found when at least 60 genes are evolving as neutral sequences (**fig. 5A**).

352   Using the same approach, in *L. gibarensis*, we analysed the observed distribution of the

353   number of LoF mutations per pseudogene (71 genes without LoF mutation, 3 with 1 mutation,

354   2 with 2 mutations), considering a number of neutral genes within a range of 5 to 76 (**fig. 5B**

355   and **supplementary fig. S7**, **Supplementary Material** online). In this case, the best fit was

356   obtained when about 15 eye genes are free to accumulate LoF mutations (**fig. 5B**). These

357   results suggested that most genes, if not all, are dispensable in the blind *L. dentata* whereas

358   only a small subset can be lost in the small-eyed *L. gibarensis*.

359

## Evidence of relaxed selection on non-synonymous mutations in cavefish eye genes

362

363  To reinforce the evidence brought by the above analyses on LoF mutations, we looked for

364  other signatures of relaxed selection using methods based on changes in ω (the ratio of the

365  mean number of nonsynonymous substitutions per nonsynonymous site to the mean number

366  of synonymous substitutions per synonymous site, also known as dn/ds). It is expected to be

367  lower than one under purifying selection, equal to one under neutral evolution, and larger than

368  one under adaptive selection. As gene divergence between *Lucifuga dentata* and *Lucifuga*

369  *gibarensis* was lower than 0.9% and lower than 0.2% between the two *Astyanax mexicanus*

370  morphs (for more details, see **supplementary folder divergence_values**, **Supplementary**

371  **Material** online), the number of nucleotide differences per gene was very low and often no

372  sequence change was observed between a cave species (or population) and the closest surface

373  species (or population) (**supplementary fig. S12**, **Supplementary Material** online).

374  Therefore, we used three sets of concatenated gene sequences (eye, circadian clock and

375  pigmentation genes) to compute ω.

376  With the phylogenetic analysis using maximum likelihood (PAML) package Version 4.9h

377  (Yang 2007), allowing a different ω along each branch, *Lucifuga dentata* had the highest ω

378  (0.409) for eye genes. For circadian clock genes, both *Astyanax mexicanus* cavefish and

379  *Lucifuga dentata* had the highest ω (0.29). For pigmentation genes, ω was similar in cave and

380  surface fishes (**fig. 6ABC**, **supplementary fig. S13**, **Supplementary Material** online).

381  Independently, we computed ω for the same sets of genes in *Sinocyclocheilus* species. For

382  each species, ohnologs were concatenated into two series of gene sequences. With eye genes,

383  ω was higher in the blind *S. anshuiensis* (0.36) than in the small eyed *S. rhinocerous* (0.32)

15

384  and the eyed *S. grahami* (0.23). With circadian clock genes, ω has higher in the blind *S.*

385  *anshuiensis* (0.38) and the small eyed *S. rhinocerous* (0.37) than in the eyed *S. grahami*

386  (0.25). With pigmentation genes ω was higher in the small eyed *S. rhinocerous* (0.32) and the

387  blind *S. anshuiensis* (0.29) than in the eyed *S. grahami* (0.25) (**supplementary fig. S13**,

388  **Supplementary Material** online). Thus ω was consistently higher in cavefishes than in

389  surface fishes, the shift being larger for eye genes than for circadian and pigmentation genes.

390  In order to further examine if the shift of ω in some cavefishes for some sets of genes revealed

391  a relaxed selection, we used another approach implemented in RELAX which computes the

392  values and distribution of three ω using a branch-site model, the convergence of the three ω

393  towards one in a lineage being a signature of relaxed selection (Wertheim, et al. 2015). The

394  magnitude of convergence depends on a parameter, k, which tends to zero as selection tends

395  to complete relaxation. RELAX detected relaxed selection on *Lucifuga dentata* eye genes

396  with an important shift toward ω = 1 (k = 0.2), and this was also true in a lesser extent in *A.*

397  *mexicanus* cavefish (k = 0.5). For pigmentation genes, the largest shift was also observed in

398  *Lucifuga dentata* (k = 0.48). No shift of distribution was observed with cavefish circadian

399  clock genes, suggesting that these genes are under strong purifying selection (**supplementary**

400  **fig. S14, fig. S15** and **fig. S16 Supplementary Material** online).

401  Finally, with the aim of finding independent evidence of relaxed selection in cavefishes, in

402  particular on *A. mexicanus* eye genes for which the number of mutations is low and thus the

403  estimation of ω was not accurate, a novel approach was developed. First, nonsynonymous

404  mutations in different lineages were inferred using the aaml program from the PAML

405  package. Then, the deleterious impact of these mutations, a score which ranges between 0 (not

406  deleterious) and 1 (very deleterious), was estimated using a machine learning method

407  implemented in MutPred2 (Pejaver, et al. 2017). The kernel density estimation (KDE) of the

408  distributions of the scores in eye, circadian clock and pigmentation genes were obtained for

409    each terminal lineages leading to surface fishes and cavefishes, as well as for computer

410    simulations of substitutions in the same sets of genes under a neutral model. Whatever the set

411    of genes, in all surface fishes, the KDE was similarly right-skewed (**fig. 7ABC**), suggesting

412    that most mutations which reached fixation have a low impact on the fitness. This was

413    confirmed by the shape of the distribution of the scores in simulations of substitutions without

414    selection (equivalent to the distribution before selection) which was very different to those of

415    surface fishes, that is almost uniform and suggesting that the most deleterious mutations had

416    been removed by selection in surface fishes. Before selection, the score distribution was

417    slightly different for the different sets of genes, probably reflecting different selective

418    constraints on the sequences belonging to these gene sets (**fig. 7ABC**, grey and black curves).

419    Noteworthy, the Transitions/Transversions (Ts/Tv) ratio used in simulations of substitutions

420    under a neutral model had no impact on the distribution of the scores (**supplementary fig.**

421    **S19, Supplementary Material** online). In cavefishes, the score distribution was very

422    variable, depending on the cavefish species and the set of genes (**fig. 7ABC**). In order to

423    refine the analysis of the score distribution in cavefishes, admixtures of different proportions

424    of substitutions picked up from two distributions, one under neutral evolution (from the

425    simulations) and the other with selection (in the lineage leading to zebrafish) were also

426    obtained to make comparisons with cavefish distributions. Pairwise comparisons of empirical

427    cumulative distributions (ECDF) were performed using the nonparametric Kolmogorov-

428    Smirnov (KS) test. The same approach was attempted using Grantham's distances (Grantham

429    1974) instead of MutPred2 scores but the contrast between the distributions of the distances

430    with and without selection was much less discriminant and not analyzed further

431    (**supplementary fig. S18, Supplementary Material** online).

432    With eye genes, for *A. mexicanus* cavefish (red curve, **fig. 7**), the distribution was not

433    statistically different from that expected if all substitutions where neutral in this lineage (KS

17

434    test, p = 0.2; **supplementary fig. S17, Supplementary Material** online), yet the best fit was

435    with a mixture distribution with 24% of substitutions from the distribution under selection

436    (**supplementary fig. S20, Supplementary Material** online). For *L. dentata* (brown curve,

437    **fig. 7**) and *L. gibarensis* (orange curve, **fig. 7**), distributions departed from the neutral

438    distribution (KS test, p = 1.4 x $10^{-5}$ and p = 4 x $10^{-6}$ respectively) (**fig. 7A** and **supplementary**

439    **fig. S17, Supplementary Material** online) and the best fit was obtained with respectively

440    34% and 60% of the substitutions from the distribution under selection (**supplementary fig.**

441    **S20, Supplementary Material** online). For all *Sinocyclocheilus* species, the score

442    distribution was different from those of surface fishes, even for the eyed *S. grahami*, most

443    likely because after the WGD purifying selection on nonsynonymous mutations was partially

444    relaxed on one or both ohnologs, but the ECDF of *S. rhinocerous* and *S. anshuiensis* were

445    more shifted toward the neutral distribution than the ECDF of *S. grahami*, suggesting that the

446    two cavefishes experienced a more neutral regime than the surface fish (**supplementary fig.**

447    **S21, Supplementary Material** online).

448    <u>With circadian genes</u>, no cavefish ECDF fitted with the expected distribution under neutral

449    evolution (**fig. 7B**). However, the ECDF of *A. mexicanus* cavefish was different from those of

450    surface fishes and the best fit was obtained with an admixture of 59% of the substitutions

451    from the distribution under selection (**supplementary fig. S20, Supplementary Material**

452    online). For *L. dentata* and *L. gibarensis*, the best fit involved respectively the admixture of

453    69% and 93% of the substitutions from the distribution under selection (**fig. 7B,**

454    **supplementary fig. S20, Supplementary Material** online). In accordance with the number

455    of pseudogenes found in *S. rhinocerous* for this set of genes, its ECDF was the closest to the

456    neutral distribution among the three *Sinocyclocheilus* species, with a best fit found with an

457    admixture of 39% of substitutions from the distribution under selection (**supplementary fig.**

458    **S21** and **fig. S22, Supplementary Material** online).

459   With pigmentation genes, no cavefish ECDF fitted with the expected distribution under

460   neutral evolution (**fig. 7C**). All cavefish distributions were very similar to the surface fish

461   distributions, in accordance with the hypothesis that very few genes belonging to this category

462   can be lost, even after cave colonization and/or genome duplication (**see also supplementary**

463   **fig. S17, fig. S20, fig. S21 and fig. S22, Supplementary Material** online).

464   In summary, the different approaches consistently suggested that different levels of relaxed

465   selection on the set of eye genes are correlated with the levels of eye degeneration in

466   cavefishes, whereas most circadian clock and pigmentation genes are under strong purifying

467   in these species.

468

469   <span style="color:purple">Dating relaxation of purifying selection on eye genes in *L. dentata*</span>

470

471   In order to conciliate results suggesting that most eye genes are dispensable and the finding

472   that selection is not totally relaxed in the *L. dentata* lineage, we postulated two successive

473   periods of evolution, one under selection followed by another under relaxed selection. Three

474   independent approaches were used to estimate when selection was relaxed in this lineage.

475   First, we used the number of eye pseudogenes. With a simple analytical model assuming a

476   LoF mutation rate equal to $0.072 \times 10^{-8}$, the highest probability of finding 19 pseudogenes

477   among 76 neutral genes was obtained for relaxed selection starting 367,779 generations ago

478   (probability > 5% in a range between 273,990 and 480,980 generations) (**fig. 8**, red curve).

479   Assuming that only 50 eye genes were free to accumulate LoF mutations, this time was

480   pushed back to 611,132 [445,950 – 813,580] generations (**fig. 8**, pale red curve). Then,

481   simulations were performed in order to take into account variations in the gene length and the

482   number of introns per gene (**supplementary Data_Supp1, Supplementary Material** online),

483   codon usage, transition/transversion ratio (r = 4.57) and effective population size ($N_e$) in a

484  range between 100 and 1,000. These simulations and the analytical model gave very similar

485  estimations, the effects of $N_e$ and per gene LoF mutation rate variation being marginal (**fig. 8**,

486  black, green and blue curves, only simulations assuming 76 neutral genes are shown).

487  Second, two dating methods were used (Li, et al. 1981; Meredith, et al. 2009), both based on

488  the hypothesis of a shift of ω from a low value to 1 after purifying selection was relaxed in a

489  lineage. We assumed a divergence time of 80 My between *Lucifuga* and *Brotula*

490  (http://www.timetree.org/). Eye genes of *Lucifuga* species and B*rotula barbata* were

491  individually aligned and alignments concatenated. With one method (Li, et al. 1981), the

492  divergence time between *Lucifuga dentata* and *Lucifuga gibarensis* was estimated equal to

493  4,110,441 years ago and the time since non-functionalization of eyes genes in *L. dentata* equal

494  to 1,486,042 years. With the other method (Meredith, et al. 2009), ω was estimated equal to

495  0.271 in the lineage leading to *L. gibarensis* and 0.502 in the lineage leading to *L. dentata*.

496  Assigning these ratios respectively to functional branches and a mixed branch, the time since

497  non-functionalization was estimated to 1,302,485 years.

498  Third, assuming that in the lineage leading to *L. dentata,* there is an admixture of 66% of the

499  mutations that accumulated under relaxed selection and 34% under selection (**supplementary**

500  **fig. S20, Supplementary Material** online), and that ω = 0.27 under selection (that is ω

501  estimated in *L. gibarensis*, see **supplementary Data_Supp3, Supplementary Material**

502  online) and ω = 1 under relaxed selection, and the divergence between *L. dentata* and *L.*

503  *gibarensis* occurred 4,110,441 years ago (estimated above), using the approach described in

504  Materials and Methods, we obtained a congruent estimation of the age of selection relaxation,

505  that is 1,413,991 years ago (**table 1, supplementary Data_Supp3, Supplementary Material**

506  online). Thus, the various methods to date relaxation of purifying selection in *L. dentata*

507  lineage converge to approximately 1.5 My or 400.000 generations, estimations that would be

508  consistent if we assume a generation time of about 4 years in *Lucifuga* cavefishes.

509

## Distribution of LoF mutations in eye genes of cavefishes *vs* fossorial mammals

511

512    An extensive study of the regression of visual protein networks in three fossorial mammals,

513    the Cape golden mole *Chrysochloris asiatica*, the naked mole-rats *Heterocephalus glaber* and

514    the star-nosed mole *Condylura cristata*, has been published (Emerling and Springer 2014).

515    From this publication, we retrieved the number of pseudogenes, their names, and the number

516    of LoF mutations per pseudogene in the three species (**fig. 9A**). In the Cape golden mole, 18

517    pseudogenes were found among 63 eye genes, while only 11 pseudogenes were found in the

518    naked-mole rat and 7 in the star-nosed mole. The distributions of LoF mutations per

519    pseudogene in these mammals and those of two blind cavefishes (*L. dentata* and *S.*

520    *anshuensis*) were compared (**fig. 9B**). Many independent LoF mutations were found in the

521    same eye genes in fossorial mammals and in cavefishes (**fig. 9A**). For each species, the

522    distribution of the number of LoF mutations per pseudogene, either taking into account only

523    shared pseudogenes between mammals and fishes or all the pseudogenes, were similar (**fig.**

524    **9B**, main graph and inset respectively). However, the distribution was sharply contrasted

525    between mammals and fishes. In fossorial mammals, most pseudogenes carried many LoF

526    mutations, up to 28 mutations in two pseudogenes of the golden mole and 54 mutations in a

527    single pseudogene of the star-nosed mole (**fig. 9**, **supplementary Data_Supp1,**

528    **Supplementary Material** online). On the contrary, in fishes, very few LoF mutations were

529    found in each pseudogene (**fig. 9**, **fig. S7**, **supplementary Data_Supp1, Supplementary**

530    **Material** online), the maximum being 5 LoF mutations in one pseudogene of *S. anshuiensis*

531    and 3 LoF mutations in a pseudogene of *L. dentata*. This comparison strongly supports the

532    hypothesis that fossorial mammals have lived in the absence of light for a much longer time

533    than cavefishes, but a smaller subset of genes has been under relaxed selection in mammals.

534

## Discussion

536

537 When selection for maintaining a protein in a functional state is relaxed, theory predicts that

538 LoF mutations in its coding and regulatory sequences can reach fixation by random genetic

539 drift (Lynch and Conery 2000; Lahti, et al. 2009). In an isolated population, among a set of

540 dispensable genes, the longer the time of neutral evolution, the higher the expected number of

541 pseudogenes. Eventually, all the genes under relaxed selection will be pseudogenized. At the

542 level of a single gene, the longer the time of neutral evolution, the higher the expected number

543 of LoF mutations. Thus, after a very long period of time of neutral evolution, all the neutrally-

544 evolving genes must carry many LoF mutations. The pace of this gene decay depends

545 essentially on the pace of appearance of LoF mutations (Li and Nei 1977). In the present

546 study, we focussed on a subset of LoF mutations that can be readily detected in genomes, that

547 is mutations generating internal STOP codons, eliminating START or STOP codons,

548 disrupting intron splice sites, and small insertions/deletions (indels) causing translation

549 frameshifts. Although this approach inevitably leads to an underestimation of the number of

550 non-functional genes, it allows comparative studies and molecular dating of selection

551 relaxation in different species. Below, we discuss the patterns of pseudogenization in different

552 sets of genes involved in vision, circadian clock and pigmentation during evolution in the dark

553 of several cavefishes. We show how pseudogenization of eye genes in *Lucifuga dentata* shed

554 new light on gene loss in relation to eye regression in cavefishes. On this basis, we refine

555 previous analyses of other cavefish genomes. At a broader phylogenetic scale, we discuss the

556 contrasted dynamic of pseudogenization in cavefishes and fossorial mammals.

557

558 Putative impact of some LoF mutations

559

560     <u>Eye genes</u>: in *L. dentata*, a frameshift was found in the alpha-crystallin, *cryaa*, whose

561 downregulation in *A. mexicanus* cavefish plays a key role in triggering lens apoptosis (Ma, et

562 al. 2014; Hinaux, et al. 2015). Another crystallin, *crybb1*, is pseudogenized in *L. dentata*.

563 Mutations in this gene cause lens opacity in humans (Mackay, et al. 2002). We also found

564 LoF mutations in two opsin receptor kinases, *grk7a* and *grk1b*. Mutations in these proteins

565 can lead to overactive opsin and photoreceptor degeneration (Feng, et al. 2017). These two

566 genes and *grk7b* have similar functions and are all expressed in cones. As these three kinases

567 may have additive effect (Osawa and Weiss 2012), we can hypothesize that absence or

568 malfunction of one of them could be compensated by the others. Such compensation could

569 explain why we found that both *grk7b* ohnologs carry LoF mutations in *S. grahami,* despite

570 this fish has large eyes showing no evidence of degeneration. Another interesting gene is

571 *gnb3b* which is pseudogenized in both *L. dentata* and *L. gibarensis* and which is linked to

572 night-blindness in humans (Vincent, et al. 2016), yet *gnb3$^{-/-}$* mice seem to have functional

573 photoreceptors. Finally, we found LoF mutations in *gcap2*, a guanylate cyclase activator, in

574 both *Lucifuga* species. This gene is associated to *retinitis pigmentosa* in humans (Sato, et al.

575 2005) but it could be compensated by overexpression of *gcap1* in rods (Makino, et al. 2012).

576 In *Astyanax mexicanus*, a deletion of 11 bp in the phosphodiesterase *pde6b,* a rod-expressed

577 gene, leads to several STOP codons in the catalytic domain (Lagman, et al. 2016). Mutations

578 in this gene were associated with night-blindness and *retinitis pigmentosa* in humans

579 (McLaughlin, et al. 1993; Gal, et al. 1994). Moreover, in mice affected by mutations in the

580 ortholog of *pde6b*, rod photoreceptors degenerate during development resulting in a total

581 absence of photoreceptors in the adult (Farber and Lolley 1974; Chang, et al. 2002).

582 Most LoF mutations were found in the subset of non-visual opsins, which makes their

583 functional impact difficult to evaluate as the functions of these genes are still poorly

584    understood. Two notable exceptions are *opn4m2* and *tmt3a,* pseudogenized in *S. anshuiensis*

585    and *L. gibarensis* respectively, and known to be non-functional and as such involved in the

586    deregulation of the circadian clock in *P. andruzzii*.

587

588    Circadian clock genes: in *S. rhinocerous,* both ohnologs of four circadian clock genes, *cry1b*,

589    *cry2a*, p*er2* and *cry-dash,* carried LoF mutations. In *S. anshuiensis*, both ohnologs of *cry-dash*

590    carried also LoF mutations which are independent from those found in *S. rhinocerous*. The

591    gene *cry-dash,* involved in photoreactivation DNA repair, is also pseudogenized in

592    *Phreatichthys andruzzii*  (Zhao, et al. 2018) as well as *per2* that could be involved in the

593    disruption of the circadian rhythm in this species (Ceinos, et al. 2018).

594

595    Pigmentation genes: both *L. dentata* (depigmented skin) and *L. gibarensis* (pigmented skin)

596    carried independent LoF mutations in *myo7ab*. While no *myo7ab-/-* mutant has been

597    analyzed, the paralog *myo7aa-/-* mutant in zebrafish showed an elevated photoreceptor death

598    but pigmentation was not affected (Wasfy, et al. 2014). Both *Lucifuga* species had

599    independently fixed LoF mutations in *smtla* which is known to increase the number of

600    leucophores at the expense of a reduced number of xanthophores in medaka (Fukamachi, et

601    al. 2009). In *L. dentata, slc2a11b* is pseudogenized and this gene codes for a protein that

602    promotes yellow pigmentation (Kimura, et al. 2014; Parichy and Spiewak 2015). Two other

603    genes, *trpm1a* and *trpm1b* are also pseudogenized in *L. dentata*. During zebrafish

604    development, *trpm1a* is expressed in the retina and melanophores whereas *trpm1b* expression

605    is restricted to the retina (Kastenhuber, et al. 2013). In human, mutations in their ortholog

606    TRPM1 lead to complete congenital stationary night blindness (Audo, et al. 2009). In *L.*

607    *gibarensis*, *pax7* which promotes xanthophore differentiation (Nord, et al. 2016) carried a LoF

608    mutation as well as *edn3b* that is known to lead to a reduction in iridophore numbers when

24

609  mutated in zebrafish (Krauss, Frohnhöfer, et al. 2014). In *Astyanax mexicanus*, two

610  pigmentation genes were found with LoF mutations: *mc1r* which carried a 2 bp deletion that

611  could be involved in pigmentation reduction in two cave populations belonging to this species

612  (Gross, et al. 2009) and *tyrp1a* which carried a 1 bp deletion. In zebrafish, morpholino-

613  induced knock-down of *tyrp1a* had no phenotypic effect (Krauss, Geiger-Rudolph, et al.

614  2014). In *S. rhinocerous* (pigmented skin) and *S. anshuiensis* (depigmented skin) both

615  ohnologs of *gch2* and *pmelb* carried independent LoF mutations. It has been shown that *gch2*

616  mutant lacked proper xanthophore pigmentation at larval stages in zebrafish but no effect

617  were reported in the adult (Parichy, et al. 2000; Pelletier, et al. 2001; Lister 2019). In the same

618  way, injection of *pmelb* morpholinos in the zebrafish had no significant effect on the number

619  of melanosome but led to a significant loss of their cylindrical shape (Burgoyne, et al. 2015).

620  Many pigmentation pseudogenes seem to be compensated by their teleost-specific duplicates

621  when lost in zebrafish, such as *tyrp1a* (Krauss, Geiger-Rudolph, et al. 2014), *pmelb*

622  (Burgoyne, et al. 2015) and *pax7b* (Nord, et al. 2016).

623

624  Contrasted decay of eye genes *vs* circadian clock and pigmentation genes

625

626  In order to study pseudogenization in relation to the regression of three traits in cavefishes, we

627  defined three categories, that are eye, circadian clock and pigmentation genes. For most

628  genes, assigning a gene to a category was straightforward, yet for some genes it was more

629  ambiguous. Most eye genes corresponded to a set of genes expressed only in eyes, however

630  fishes also express several non-visual opsins genes that we assigned to this category on the

631  basis of their homology to visual opsins. Genes known for being involved in the circadian

632  clock were assigned to a second set of genes. Noteworthy, some non-visual opsins are

633  involved in this process. Pigmentation genes comprised a large set of genes involved in

634    several processes from pigment cell differentiation to pigment synthesis. Our *a priori*

635    hypothesis was that eye genes should be more prone to degenerate in blind fishes as there are

636    only expressed in eyes or involved in light sensing in other tissues, whereas many circadian

637    clock and pigmentation genes may be maintained as their expression is not restricted to

638    regressed structures and have pleiotropic roles. Indeed, while many pseudogenes were

639    identified among eye genes of some cavefishes, a much smaller proportion of pseudogenes

640    were found among circadian clock and pigmentation genes. In addition, several cases of

641    parallel fixation of LoF mutations in different species among a small subset of genes

642    suggested that only few genes involved in the circadian clock and pigmentation can be lost in

643    cavefishes.

644

645    Molecular evidence of circadian clock disruption in several cavefishes

646

647    No LoF mutations were found in the set of circadian clock genes of both *Lucifuga* species.

648    However, *tmt3a*, a non-visual opsin is pseudogenized in *L. gibarensis* and the loss of this gene

649    is involved in the disruption of the circadian rhythm in the cavefish *Phreatichthys andruzzii*.

650    Whereas the survey of LoF mutations did not allow to find evidence of circadian clock loss in

651    *L. dentata*, it is probably the case in *L. gibarensis*. Selection on circadian clock genes is also

652    supported by the analysis of non-synonymous mutations which suggested no higher

653    deleterious mutation accumulation in these species when compared with surface fishes.

654    As expected, no LoF mutations in both ohnologs of circadian clock genes and non-visual

655    ospin genes was found in *S. grahami* which is a surface fish. Unexpectedly, the small-eyed *S.*

656    *rhinocerous* has accumulated more circadian clock pseudogenes (*per2*, *cry-dash*, *cry1b*,

657    *cry2a*) than the blind *S. anshuiensis* (*cry-dash*), suggesting that the level of eye regression

658    could be loosely correlated with the level of circadian clock disruption. Moreover, as several

659     independent LoF mutations were found in a small number of circadian clock genes, some of

660     them already known to be involved in the circadian clock disruption in other species, it

661     suggests that pseudogenization of a small subset of genes can be involved in this process, in

662     particular those belonging to cryptochromes and period families which are light-inducible

663     genes.

664

665     ## A small subset of pigmentation pseudogenes

666

667     A similar trend was observed among pigmentation genes: independent LoF mutations were

668     found in *myo7ab* and *smtla* of *L. dentata* and *L. gibarensis* and both ohnologs of *gch2* and

669     *pmelb* carried independent LoF mutations in *S. anshuiensis* and *S. rhinocerous*. Recurrent

670     pseudogenization of the same genes suggests that a very small subset of pigmentation genes

671     can be lost, and that these genes might be those which have no or few pleiotropic effects.

672     Indeed, many pigmentation genes code for transcription factors or signaling molecules

673     involved in neural crest-derived, pigment cell differentiation, that are repeatedly used at

674     different times and places during development (Betancur, et al. 2010).

675

676     ## Many eye pseudogenes in the ancient diploid cavefish *Lucifuga dentata*

677

678     Before the present study, there was no evidence on the possibility of pseudogenization of

679     many eye genes in blind cavefishes. In *L. dentata*, we found up to 25% of eye genes carrying

680     LoF mutations. Moreover, the distribution of LoF among genes is consistent with neutral

681     evolution of a large proportion of, if not all, eye genes in this species. On the other hand, in *L.*

682     *gibarensis* which has small but functional eyes, most eye genes seem under selection but the

683     partial degeneration of the visual system is correlated with the loss of several genes well

684    conserved in eyed fishes. These data allowed us to propose a two-step scenario for the release

685    of selection pressure on eye genes in this genus. The common ancestor of *L. dentata* and *L.*

686    *gibarensis* was an eyed cavefish that had accumulated a small number of pseudogenes in

687    relation to life in darkness, but none among eye specific genes. In *L. gibarensis*, most eye

688    specific genes have been under purifying selection whereas it has been relaxed in *L. dentata*.

689    Interestingly, the population of *A. mexicanus* from the Pachón cave which is very recent but in

690    which cavefish have highly degenerated eyes, only one eye gene carried a LoF mutation. The

691    lack of correlation between the degree of eye regression and the number of eye pseudogenes

692    underscores the fact that the extent of visual regression should not be taken as a proxy of the

693    evolutionary age of cavefish populations or species.

694

695    Dating blindness in *L. dentata*

696

697    Dating changes in selective constraints on traits and genes after cave settlement is a difficult

698    task. Several closely-related methods have been proposed to estimate when a change of

699    selective regime occurred on one gene in one lineage, that is when ω shifted from a value

700    lower than one (a signature of purifying selection) to one (a signature of neutral evolution)

701    (Li, et al. 1981; Miyata and Yasunaga 1981; Meredith, et al. 2009; Zhao, et al. 2010;

702    Wertheim, et al. 2015). With two different methods (Li, et al. 1981; Meredith, et al. 2009), we

703    estimated that the time since selection was released on the eye genes of *Lucifuga dentata* is

704    between 1.3 Mya and 1.5 Mya. Taking into account that 19 pseudogenes were found among

705    76 eye genes that may be dispensable for a blind fish, and assuming a LoF mutation rate equal

706    to $0.072 \times 10^{-8}$ per site per generation, we estimated the time since *L. dentata* settled in caves

707    about 380,000 generations ago. The generation time is unknown for this fish, and translating

708    the number of generations into years is difficult. However, assuming that the generation time

28

709    is about four years, which is realistic if we consider that they could reproduce during about

710    ten years, the above independent estimations of relaxed selection would be coherent.

711    Moreover, using the distribution of the MutPred2 scores, we obtained another and very close

712    estimation (1.4 Ma). Our results suggest that *L. dentata* and *L. gibarensis* could have diverged

713    more than 4 million years ago. The common ancestor of these species could have had well

714    developed eyes that slightly regressed in one lineage (*L. gibarensis*) but much more in the

715    other (*L. dentata*) after a long period without degeneration; or else, the ancestor could have

716    had small eyes like *L. gibarensis* which after a long stasis completely degenerated in the

717    lineage leading to *L. dentata* but remained almost unchanged in the lineage leading to *L.*

718    *gibarensis*.

719    Thus, the magnitude of eye degeneration that is often used as a proxy of the age of cave

720    species because it is assumed that eyes degenerate gradually and continually in such

721    environment is likely often misleading. A refined analysis of fish ecology is necessary to

722    better understand the pace and the level of eye degeneration. Indeed, caves are often described

723    as repetitions of the same environment, that is highly isolated and totally dark. However,

724    some cavefishes such as *Lucifuga gibarensis* and closely-related small eyed species can be

725    found in caves that are partially lighted, or sink holes in the sea. Such a complex environment

726    could be the reason for the maintenance of small yet functional eyes in these species, like in

727    fossorial mammals.

728

729    Pattern of LoF mutations in recent tetraploids with different level of

730    troglomorphy: the case of *Sinocyclocheilus*

731

732    The genus *Sinocyclocheilus*, which is endemic to southwestern karst areas in China, is the

733    largest cavefish genus known to date (Xiao, et al. 2005). LoF mutations were found in several

734    genes of three species, one species (*S. anshuinensis*) being blind and depigmented, another

735    species (*S. rhinocerous*) having small eyes and being pigmented, and the last one (*S. grahami*)

736    showing no such troglomorphic traits (Yang, et al. 2016). These species share a WGD with

737    other cyprinids such as the common carp *Cyprinus carpio* (David, et al. 2003; Yuan, et al.

738    2010) which could explain why even the surface fish carry many LoF mutations in eye,

739    circadian clock and pigmentation genes (Yang, et al. 2016). However, no thorough

740    comparisons were performed. Our results are consistent with a rapid radiation within this

741    genus (Xiao, et al. 2005) as only few LoF mutations were found in internal branches of their

742    phylogenetic tree. The divergence between *Cyprinus carpio* and the *Sinocyclocheilus* species

743    may have occurred soon after the WGD as only two shared LoF mutations were found. The

744    number of eye pseudogenes in the blind *S. anshuiensis* is much higher than in the small-eyed

745    *S. rhinocerous* and the eyed *S. grahami*, a result supporting the cumulative effect of

746    tetraploidy and cave settlement on the rate of accumulation of LoF mutations. As most genes

747    are present twice, a gene function is lost if, and only if, at least one LoF mutation is present in

748    each ohnolog. With this criterion, seven genes were lost in *S. anshuiensis*, but only one gene

749    in *S. rhinocerous* and *S. grahami*. Selective pressure was relaxed on one copy of these genes

750    after the WGD, but a complete relaxation occurred only after cave settlement in *S.*

751    *anshuiensis*. Among the genes for which both ohnologs are mutated in *S. anshuiensis*, the

752    mutations in *pde6c* could have a role in photoreceptors degeneration, as suggested by a study

753    of zebrafish mutants (Stearns, et al. 2007). *Sinocyclocheilus rhinocerous* lost the two

754    functional copies of *gcap1* and it has been shown that two missense mutations in this gene

755    lead to significant disruptions in photoreceptors and retinal pigment epithelium, together with

756    atrophies of retinal vessels and choriocapillaris in zebrafish (Chen, et al. 2017). However,

757    knockout of *gcap1* in mice showed that its absence does not change expression level of other

30

758    phototransduction proteins thanks to a compensation by *gcap2*. Nevertheless, the knock-down

759    leads to a delayed recovery after light exposure (Makino, et al. 2012).

760    Analyses with RELAX and the estimation of the admixture of MutPred2 score distributions

761    that best fit with the observed score distribution also suggest that purifying selection on eye

762    genes is much higher in *S. grahami* than in *S. anshuiensis* and *S. rhinocerous*, much lower on

763    circadian genes of *S. rhinocerous* but high on pigmentation genes in these three species.

764    These results are congruent with the level of pseudogenization observed for the three gene

765    sets in the three species.

766

767    Very few pseudogenes in the recent settler *Astyanax mexicanus*

768

769    In the reference genome of *Astyanax mexicanus* cavefish, a LoF has been found in the eye

770    gene *pde6b*. This mutation went unnoticed in previous studies but may well contribute to

771    retinal degeneration. No LoF were found in clock genes. Among pigmentation genes, a 1 bp

772    deletion was found in *tyrp1a* and a 2 bp deletion in *mc1r*. The latter mutation has been

773    associated with the brown phenotype of some populations (Gross, et al. 2009) but the finding

774    of a close and functional tandem duplicate suggest that it actually may not be the cause of this

775    phenotype (Gross, et al. 2017). Overall, these results are in accordance with a very recent

776    settlement of *Astyanax* cavefish (Fumey, et al. 2018) that did not allow the fixation of many

777    eye pseudogenes despite the lack of purifying selection on most, if not all, eye genes. The

778    extreme eye degeneration with only one LoF in *Astyanax* cavefish eye genes further questions

779    the nature of the developmental mechanisms involved in eye loss in this species, the pace of

780    eye degeneration and the correlation of eye degeneration with gene decay.

781

782    Contrasted dynamics of pseudogenization in fossorial mammals and cavefishes

783

784    The genomes of three independently-evolved fossorial mammals have previously allowed an

785    extensive study of LoF mutations in genes coding for proteins involved in retinal networks

786    (Emerling and Springer 2014). These animals have functional eyes, but star-nosed moles

787    often leave their burrows and have the greatest exposure to light whereas naked mole-rats and

788    Cape golden moles are entirely subterranean. In addition, the eyes of Cape golden moles are

789    subcutaneous. More pseudogenes were found in the Cape golden mole than in the naked-rat

790    genome and the lowest number of pseudogenes was found in the star-nosed mole genome,

791    suggesting that the decrease in retinal exposure to light allowed the decay of more eye genes.

792    The most striking difference between cavefishes and fossorial mammals is that pseudogenes

793    of cavefishes accumulated only one or a couple of LoF mutations per pseudogene whereas

794    some pseudogenes of fossorial mammals carried a large number of LoF mutations. This

795    difference in molecular decay strongly suggests that the fossorial mammals adapted to the

796    subterranean environment a long time ago whereas colonisation of the dark environment is

797    much more recent in the case of the cavefishes.

798

799    Conclusion

800

801    Our analyses suggest that blind cavefishes examined so far are not very ancient. They all lost

802    their eyes during the Pleistocene, the oldest during early Pleistocene and the most recent

803    during the late Pleistocene or even later in the Holocene. The sequencing of a large number of

804    blind cavefish genomes will be necessary to identify the whole set of eye genes that are

805    dispensable in the dark, when eyes are highly degenerated. Moreover, finding a blind cavefish

806    genome in which most eye genes are pseudogenized and carry many LoF mutations would

807    refute our current working hypothesis that blind cavefishes cannot thrive for a very long time

808    in cave ecosystems.

809

## **Materials and Methods**

811

### Assembly of *L. dentata and L. gibarensis* draft genomes

813

814    The sequenced *L. dentata* specimen was a female, blind and depigmented. All the fish

815    belonging to this species are blind whereas their pigmentation is highly variable (Garcia-

816    Machado, unpublished data). The sequenced *L. gibarensis* specimen was a male, had small

817    eyes and was pigmented. All the fish belonging to this species have small eyes whereas their

818    pigmentation is also highly variable (García-Machado, et al. 2011). DNA was extracted using

819    a protocol already described elsewhere (García-Machado, et al. 2011). For *L. dentata*, paired-

820    end libraries were prepared with different insert sizes: 200 bp, 400 bp and 750 bp. A mate-

821    pair library was also prepared with insert size in the range 3-5 kb. For *L. gibarensis*, only one

822    mate-pair library was prepared, which had inserts size between 3 kb and 10 kb. *Lucifuga*

823    *dentata* libraries were sequenced on an Illumina HiSeq 2000 sequencer whereas *L. gibarensis*

824    library was sequenced on an Illumina NextSeq sequencer. After cleaning steps (adaptors

825    trimming and quality trimming), *L. dentata* assembly of a draft genome was performed using

826    Minia (Chikhi and Rizk 2013) on all data, resulting in 662,154 contigs. After assembling, and

827    as Minia doesn't use the paired-end information, scaffolding steps were performed using

828    SSPACE (Boetzer, et al. 2011) on one library at a time in ascending order of insert size. The

829    number of scaffolds decreased from 662,154 to 161,599 with the first library (insert size of

830    200 bp), to finish with 48,241 scaffolds with the mate pair library. This result was corrected

831    by REAPR (Hunt, et al. 2013) to obtain 52,944 scaffolds. The remaining gaps were filled by

832    GapCloser (Luo, et al. 2012).

833    The quality and completeness of the draft genome of *L. dentata* were assessed by remapping

834    paired-end reads to the assembly using BWA v0.7.11 (Li and Durbin 2009) and BUSCO

835    (Kriventseva, et al. 2015) with the Actinopterygii dataset comprising a total of 4,584

836    conserved genes. The latter analysis was performed also on published draft genomes of three

837    other Ophidiiformes (*Brotula barbata*, *Carapus acus* and *Lamprogrammus exutus*).

838    Sequences from *L. gibarensis* were mapped on the genome of *L. dentata* using BWA v0.7.11.

839

840    Assembly of *Lucifuga dentata* transcriptome

841

842    Gonads, gills, heart and brain were dissected and stored in RNA-Later (Ambion). Total RNA

843    isolation (using Trizol) lead to yields of 870 ng/µl in gonads, 750 ng/µl in gills, 240 ng/µl in

844    heart, 390 ng/µl in brain. ARN from gonads, gills and heart were mixed in equal proportions

845    to construct the first library. ARN from the brain was used to construct the second library.

846    For library preparation, polyA + RNA were extracted, fragmented, and directional libraries

847    were prepared using the Small RNA Sample Prep Kit (Illumina). Both libraries were

848    sequenced on an Illumina NextSeq500, on a Paired-end 2x150 bp run, using the High Output

849    Kit 300 cycles sequencing kit. After cleaning steps (adaptors trimming and quality trimming),

850    a *de novo* transcriptome assembly was obtained using Trinity and a quality assessment was

851    realized following Trinity recommendations

852    (https://github.com/trinityrnaseq/trinityrnaseq/wiki/Transcriptome-Assembly-Quality-

853    Assessment).

854

855    Annotation of *Lucifuga dentata* draft genome

856

857     First, repetitive elements were identified using RepeatMasker v4.0.7 (Smit, et al. 2013), Dust

858     (Morgulis, et al. 2006) and TRF v4.09 (Benson 1999). A species specific *de novo* repeat

859     library was built with RepeatModeler v1.0.11 (Smit and Hubley 2008) and repeated regions

860     were located using RepeatMasker with the *de novo* and *Danio rerio* libraries. Bedtools

861     v2.26.0 (Quinlan and Hall 2010) were used to merge repeated regions identified with the three

862     tools and to soft masked the genome. Then, MAKER3 genome annotation pipeline v3.01.02-

863     beta (Holt and Yandell 2011) combined annotations and evidence from three approaches:

864     similarity with fish proteins, assembled transcripts and *de novo* gene predictions. Protein

865     sequences from 11 other fish species (*Astyanax mexicanus*, *Danio rerio*, *Gadus morhua*,

866     *Gasterosteus aculeatus*, *Lepisosteus oculatus*, *Oreochromis niloticus*, *Oryzias latipes*,

867     *Poecilia formosa*, *Takifugu rubripes*, *Dichotomyctere nigroviridis*, *Xiphophorus maculatus*)

868     found in Ensembl were aligned to the masked genome using Exonerate v2.4 (Slater and

869     Birney 2005). RNA-Seq reads were mapped to the genome assembly using STAR v2.5.1b

870     (Dobin, et al. 2013) with outWigType and outWigStrand options to output signal wiggle files.

871     Cufflinks v2.2.1 (Trapnell, et al. 2010) was used to assemble the transcripts which were used

872     as RNA-seq evidence. A *de novo* gene model was built using Braker v2.0.4 (Hoff, et al. 2016)

873     with wiggle files provided by STAR as hints file for GeneMark and Augustus trainings. The

874     best supported transcript for each gene was chosen using the quality metric called Annotation

875     Edit Distance (AED) (Eilbeck, et al. 2009). The annotation completeness of coding genes was

876     assessed by BUSCO using the Actinopterygii gene set. Homology to uniprot database was

877     used to infer functions of predicted genes with Blastp and an e-value cutoff of $1e^{-6}$.

878     Interproscan 5.35 (Jones, et al. 2014) was used to detect proteins with known functional

879     domains.

880

## Analysis of repeated elements

881

882

883  The *de novo* library of repeated elements was refined with the following procedure: removal

884  of short (<80 bp) consensus repeats; reannotation of satellite sequences as well as of putative

885  DNA or LTR transposable elements (TEs) by aligning each consensus against itself (this

886  procedure allows to visualize internal repeats); Blastn (Altschul, et al. 1990) of the library

887  against itself and removal of redundant TEs; Blastx (Altschul, et al. 1990) of « Unknown »

888  repeats against the NCBI protein database and removal of multigene families erroneously

889  identified as putative TEs; reannotation of putative SINEs according to the SINE-scan

890  program (Mao and Wang 2017). Finally, the library was manually curated: consensus

891  sequences were compared to an in-house library of transposable element proteins using

892  BlastX. Matching Unknown elements were renamed according to their hits against this

893  library. Consensus sequences showing incongruent annotations between RepeatModeler

894  automatic classification and our manual annotation were further submitted to Censor

895  (Kohany, et al. 2006). This TE library was used as repeat database for a RepeatMasker search

896  in the genome (Smit, et al. 2013). Overlaps in RepeatMasker output were discarded by

897  selecting highest scoring elements. Repeat fragments closer than 20 bp and having the same

898  name were merged. The Landscape was reconstructed from RepeatMasker align output using

899  the calcDivergenceFromAlign.pl and createRepeatLandscape.pl utilities of the RepeatMasker

900  suite.

901

## Identification of eye, circadian and pigmentation genes

903

904  The set of eye genes included all opsins, visual opsins that are expressed in eye photoreceptor

905  cells (cone and rods) but also non-visual opsins that are expressed in a wide variety of tissues.

36

906 It also comprised eye specific crystallin genes. Crystallin genes code for several families of

907 proteins that are implicated in the transparency of the lens and fine tuning its refraction index,

908 but can also have other functions, not well known for many of them (Thanos, et al. 2014).

909 Expression patterns in zebrafish reported in ZFIN database (https://zfin.org/) and in *A.*

910 *mexicanus* (Hinaux, et al. 2015) were used to identify and select a subset of eye specific

911 crystallins. Noteworthy, *crygm2* paralogs were excluded from the analysis because many

912 copies (more than 50 copies in *A. mexicanus*) were found as in other fish genomes most likely

913 allowing relaxed selection on some copies independently to relaxed selection due to

914 environmental shift. The set of eye genes also included genes coding for proteins involved in

915 the phototransduction cascade: RPE65, Arrestins, Recoverins, Transducins, PDE6, CNGA3

916 and CNGB3, GCAPs, zGCs, and GRKs. These genes code for a highly heterogeneous set of

917 proteins with regard to their structure and functions (Imanishi, et al. 2002; Wada, et al. 2006;

918 Schonthaler, et al. 2007; Matveev, et al. 2008; Nishiwaki, et al. 2008; Rätscho, et al. 2009;

919 Renninger, et al. 2011; Fries, et al. 2013; Lagman, et al. 2015; Zang, et al. 2015; Lagman, et

920 al. 2016). Only genes whose expression was restricted to the retina and/or the pineal complex

921 were retained. Sets of circadian clock and pigmentation genes were defined on the basis of

922 gene lists established in previous studies (Li, et al. 2013; Lorin, et al. 2018). The set of

923 circadian genes was completed with *ck1δa* and *ck1δb* genes which are specific kinases of *cry*

924 and *per* genes (Takahashi, et al. 2008) and *aanat1* and *aanat2* genes whose expression are

925 regulated by the circadian clock in zebrafish (Vatine, et al. 2011). The sequences of visual and

926 non-visual opsins of zebrafish were retrieved from (Davies, et al. 2015). Other eye genes,

927 circadian and pigmentation genes of zebrafish were retrieved from GenBank.

928 Series of blastn and tblastx (Altschul, et al. 1990) with zebrafish sequences were performed

929 against *A. mexicanus* surface and Pachón cave genomes (GCF_000372685.2 and

930 GCF_000372685.1 respectively), *S. grahami*, *S. rhinocerous*, *S. anshuiensis*, *P. nattereri*, *B.*

37

931 *barbata*, *C. acus* and *L. exutus* genomes (GCF_001515645.1, GCF_001515625.1,

932 GCF_001515605.1, GCF_001682695.1, GCA_900303265.1, GCA_900312935.1 and

933 GCA_900312555.1 respectively), and *L. dentata* and *L. gibarensis* genomes (this study).

934 Matching regions were extracted using samtools (Li 2011) and coding DNA sequences (CDS)

935 were predicted using Exonerate with protein sequences of zebrafish (Slater and Birney 2005).

936

## Phylogenetic analyses

937

938

939 Orthologous and paralogous relationships between genes were inferred through phylogenetic

940 analyses. First, coding sequences were aligned using MUSCLE (Edgar 2004), after having

941 taken into account indels (*i.e.* adding N where nucleotides were missing or removing

942 additional nucleotides). For each alignment, DNA sequences were translated into protein

943 sequences and a maximum likelihood phylogenetic tree was inferred using IQ-TREE

944 (Nguyen, et al. 2015) with the optimal model found by ModelFinder (Kalyaanamoorthy, et al.

945 2017) and the robustness of the nodes was evaluated with 1,000 ultrafast bootstraps (Hoang,

946 et al. 2018). The trees were rooted and visualized using iTOL (Letunic and Bork 2006).

947 Phylogenetic trees and IQ-TREE files can be found in **supplementary folder phylogenies,**

948 **Supplementary Material** online.

949

## Identification of LoF mutations

950

951

952 We classified CDS in three classes: 1) complete, 2) pseudogene (characterized by the

953 presence of at least one among the following mutations: an internal STOP codon, an indel

954 leading to a frameshift, the loss of the initiation codon, the loss of the STOP codon, a

955 mutation in a splice site of an intron), 3) incomplete. Incomplete genes can be artifacts of

38

956   different origins such as missing data, assembly errors (Florea, et al. 2011) and gene

957   prediction errors due to sequence divergence. Nonetheless, they can be real, resulting from

958   large genomic deletions. In the case of the *A. mexicanus* cavefish genome, using PCR, we

959   could check that about 85% of the incomplete genes were assembly errors (data not shown)

960   and they were not further analyzed. Given the low quality of the *A. mexicanus* cavefish

961   genome assembly compared to the surface one and in order to get good gene sequences, cave

962   reads were retrieved and mapped onto the surface genome using the NCBI remapping service.

963   This approach allowed the identification of an opsin gene repertoire (36 genes) slightly larger

964   than the one recently published (33 genes) using only the cavefish genome (Simon, et al.

965   2019). Similarly, *Lucifuga gibarensis* reads were mapped on the *Lucifuga dentata* genome.

966   Orthologous genes from a cod (*Gadus morhua*), a medaka (*Oryzias latipes*), a platyfish

967   (*Xiphophorus maculatus*), a stickleback (*Gasterosteus aculeatus*), a pufferfish

968   (*Dichotomyctere nigroviridis*), a tilapia (*Oreochromis niloticus*) and a spotted gar

969   (*Lepisosteus oculatus*) were downloaded from Ensembl (Ensembl IDs can be found in

970   **supplementary Data_SuppS2, Supplementary Material** online). For these fishes, visual

971   opsin sequences were retrieved from an extensive study at the scale of ray-finned fishes (Lin,

972   et al. 2017).

973

974   Testing randomness of LoF mutation locations along the genes

975

976   In order to evaluate whether LoF mutations were randomly distributed or clustered along the

977   genes, we used a method initially designed for estimating the randomness of intron insertions

978   (Lynch and Kewalramani 2003). We computed the effective number of gene segments

979   defined by: $n_s = 1/\sum_{i=1}^{n} s_i^2$, with *n* being the number of segments of genes separated by *n*-1

980   LoF mutations and $s_i$ being the length of the *i*th segment. As LoF mutations are found in

981    several genes with different lengths, the position of each LoF mutation was normalized by

982    dividing by the length of the coding sequence, the sum of $s_i$ was thus equal to 1 for each gene.

983    The most extreme case of LoF dispersion is the one in which all segments are of the same

984    length ($1/n$), *i.e.* the LoF mutation are regularly spaced out, yielding $n_s = n$. On the other hand,

985    if all LoF are clustered at one end of the genes, one segment approaches length 1.0, while all

986    others approach 0.0, yielding $n_s = 1$. In order to obtain the distribution of the values of $n_s$

987    under the null model of fixation of LoF at random positions, 100,000 simulations of random

988    distribution of the observed number of LoF mutations along a gene of length 1.0 were

989    performed.

990

991    Estimation of the number of eye genes under relaxed selection in *Lucifuga* spp.

992    using the distribution of LoF mutations per gene

993

994    In order to estimate the number genes under relaxed selection (*V*) in a sample of *a priori*

995    useless eye genes (*T*) in *L. dentata* and *L gibarensis*, we compared the observed distribution

996    of LoF mutations per eye gene with the expected distribution, taking into account that only a

997    fraction (*V*) of these genes are under relaxed selection and can accumulate LoF mutations and

998    that *T - V* genes are under selection and cannot carry LoF mutations. Assuming that a LoF

999    mutation has a probability $1/V$ to appear in a gene among *V* genes under relaxed selection, the

1000   probability that a gene contains *X* LoF mutations can be computed as follows:

1001
$$p(X = 0) = \frac{V}{T}\left(1 - \frac{1}{V}\right)^m + \frac{T-V}{T} \qquad if\ i = 0$$

1002
$$p(X = i) = \frac{V}{T}\frac{m!}{i!(m-i)!}\left(\frac{1}{V}\right)^i\left(1 - \frac{1}{V}\right)^{m-i} \qquad if\ i \neq 0$$

1003   where *m* is the total number of LoF mutations.

1004    In order to take into account that eye genes do not have the same length and the same number

1005    of introns and thus mutations do not have the same probability of occurring in each gene (they

1006    are more likely in a gene with several large exons and several introns than in a gene with only

1007    one short exon), we ran 10,000 simulations of the distribution of $m$ mutations in a random

1008    sample of $V$ genes taken at random among $T$ eye genes, and taking into account the length and

1009    the number of introns in each gene to estimate its relative mutation rate. The distributions of

1010    the number of LoF mutations per gene in *L. dentata* and *L. gibarensis* were compared with

1011    expected distributions obtained with the two methods described above and for different values

1012    of $V$.

1013

1014    Sequence divergence and evidence of relaxed selection in cavefishes

1015

1016    For diploid species, genes belonging to the same gene set (eye, circadian clock or

1017    pigmentation) were concatenated. In order to analyze the genes of the tetraploid

1018    *Sinocyclocheilus* species, another alignment was produced in which each ohnolog of a given

1019    gene was concatenated with one ohnolog taken at random of the other genes, leading to two

1020    sets of concatenated genes for each species. With both alignments of concatenated sequences,

1021    maximum likelihood estimates of ω were obtained using the program codeml from the PAML

1022    package (Yang 2007) with a free-ratio model allowing a different ratio for each branch

1023    (**supplementary fig. S13, Supplementary Material** online).

1024

1025    Another approach used for detecting relaxed selection was based on analyses with the

1026    program RELAX (Wertheim, et al. 2015), assigning surface fishes as reference and excluding

1027    the small eyed fish *Lucifuga gibarensis*, the eyeless fishes *Lucifuga dentata* and *Astyanax*

1028    *mexicanus* CF. Each cavefish was independently assigned as the test branch. The value of the

41

1029    parameter k which is <1 if selection is relaxed and >1 if selection is intensified was

1030    considered as evidence of a change in the selective regime (**supplementary fig. S14**, **fig. S15**

1031    and **fig. S16, Supplementary Material** online).

1032

1033    ## Inferring the deleterious impact of amino acid variants with MutPred2

1034

1035    Maximum likelihood inference of amino acids substitutions were performed using the

1036    program aaml from the PAML package (Yang 2007). For each amino acid substitution,

1037    MutPred2 scores (Pejaver, et al. 2017) and Grantham's distances (Grantham 1974) were

1038    computed to estimate the deleterious impact of the substitutions.

1039    In order to compare the distribution of scores (or distances) for a set of genes and along a

1040    branch with the distribution expected under relaxed selection, simulations of random

1041    substitutions were generated in these genes, taking into account the length of the coding

1042    sequence of each gene and the transition/transversion ratio

1043    (https://github.com/MaximePolicarpo/Molecular-decay-of-light-processing-genes-in-

1044    cavefishes/blob/master/Neutral_evolution_for_mutpred.py). MutPred2 output files can be

1045    found **supplementary folder MutPred2_results, Supplementary Material** online).

1046

1047    ## Dating relaxation of selection with the number of eye pseudogenes in *L. dentata*

1048

1049    In absence of selection, the probability of fixation of a LoF mutation, initially absent in a

1050    population, is:

1051    $$p(1,0,t) = 1 - e^{-\mu_{LoF}t} \qquad if\ N_e \ll 1/\mu_{LoF}$$

1052    where $\mu_{LoF}$ is the LoF mutation rate, $N_e$ is the effective population size and $t$ is the number of

1053    generations (Li and Nei 1977).

1054    Thus, if $\mu_{LoF}$ is identical for a set of genes, the probability that D among T genes have fixed a

1055    LoF after time $t$ is:

$$p(X = D) = \frac{T!}{D!\,(T-D)!}(1 - e^{-\mu_{LoF}t})^D\,(e^{-\mu_{LoF}t})^{T-D}$$

1056    The derivative of this function with respect to $t$ allows to find for which value of $t$ the

1057    probability $p(X = D)$ is maximal:

$$t = \frac{1}{\mu_{LoF}}\,ln\left(\frac{T}{T-D}\right)$$

1058    For a given set of genes, the rate of LoF mutation was computed as follows:

1059    *i*) The genetic code implies that among 549 (61 x 9) mutations in sense codons, 23 lead to a

1060    STOP codon, that is ~ 4% if the frequency of each codon is 1/61 and transitions are as

1061    frequent as transversions. As among those 23 mutations, 5 are transitions and 18 are

1062    transversions, the transition/transversion ratio (*r*) can be taken into account to estimate more

1063    accurately the fraction of mutation leading to a STOP codon $f = (5r + 18)/(183r + 366)$.

1064    Using a R script, the estimation of *f* was further refined by taking into account codon

1065    frequencies (frequency_new_stop.py). For eye genes, taking into account estimations of r in

1066    *Lucifuga* spp. and *Sinocyclocheilus* spp. (4.57 and 1.95 respectively) and the codon

1067    frequencies of their eye gene sequences, *f* was estimated equal to 0.031 and 0.037 respectively

1068    in these groups of species. Moreover, taking into account that 13 internal STOP codons were

1069    found in *Lucifuga* spp. and 47 in *Sinocyclocheilus* spp., we estimated a weighted mean *f* =

1070    0.036 for the whole eye gene dataset. Applying the same approach, we found *f* = 0.038 for the

1071    circadian clock genes and *f* = 0.036 for pigmentation genes (Details in **Data_Supp1,**

1072    **Supplementary Material** online). For the three datasets taken together, weighting by the

1073    length of the concatenated genes in each dataset, we estimated a global mean *f* = 0.036. For a

1074    set of coding sequences of length *l* (sum of the CDS lengths), the rate of mutation to a STOP

1075    codon $\mu_{STOP} = f\mu l$, where $\mu$ is the nucleotide mutation rate / site.

43

1076     *ii*) The rate of indels leading to frameshifts (*i.e.* indel length modulo $3 \neq 0$) relative to the rate

1077     of new STOP codons is $\frac{n_f}{n_S}$, where $n_f$ and $n_S$ are the numbers of indels leading to frameshifts

1078     and new STOP codons respectively. The rate of frameshifts is $\mu_{frameshift} = \frac{n_f}{n_S}\mu_{STOP}$.

1079     *iii*) The rate of splice site mutations is $4n_i\mu$ (where $n_i$ is the number of introns in the set of

1080     genes).

1081     *iv*) The rate of START codon loss is $3n_g\mu$ (where $n_g$ is the number of genes).

1082     *v*) The rate of STOP codon loss is $\frac{23}{27}3n_g$ (where $\frac{4}{27}$ is the proportion of mutations in a STOP

1083     codon which leads to another STOP codon).

1084     Globally, for the set of genes, the LoF mutation rate is

$$\mu_G = \left[\left(1 + \frac{n_f}{n_S}\right)fl + 4n_i + \frac{23}{27}3n_g\right]\mu$$

1085     if all genes have the same CDS length and the same number of introns.

1086     The rate of LoF mutations per gene is $\mu_{LoF} = \frac{\mu_G}{n_g}$

1087     In order to assess the effect of the high variability of gene length and intron number observed

1088     in eye genes on pseudogene accumulation through time, a program was written to simulate

1089     decay of this set of genes through accumulation of STOP codons, frameshifts, splice site

1090     mutations, initiation and STOP codon losses, taking into account the length and the number of

1091     introns in each gene. At each generation and for each gene, the probability that a new LoF

1092     appears in one ancestral and functional allele at frequency $q$ in a population of size $N_e$ is:

1093     $2N_eq\mu_{LoF}$. When a new LoF mutation appears its frequency is $\frac{1}{2N_e}$ and the total frequency of

1094     LoF mutations is $p + 1/2N_e$, where $p = 1 - q$. We assumed random mating, no selection

1095     and no migration and a constant population size. Genetic drift between two generations was

1096     simulated taking into account the new allele frequencies if a mutation occurred, and $2N_e$ (the

1097     number of alleles sampled to generate the next generation).

1098    The simulation program was written in Python

1099    (https://github.com/MaximePolicarpo/Molecular-decay-of-light-processing-genes-in-

1100    cavefishes/blob/master/SimulationScript.py).

1101

1102    Other methods for dating selection relaxation on eye genes of *L. dentata*

1103

1104    Eye genes of the two Cuban cave brotulas (*L. dentata* and *L. gibarensis*) and an outgroup

1105    (*Brotula barbata*) were concatenated and aligned. We supposed that eye genes have been

1106    under selection along the branches of the phylogenetic tree, except in the lineage leading to *L.*

1107    *dentata* which is a mixed branch (with a period of time under selection followed by a period

1108    of time under relaxed selection). The time since selection was relaxed was estimated using

1109    two slightly different methods both relying on a shift of the nonsynonymous substitution rate

1110    after relaxed selection (Li, et al. 1981; Meredith, et al. 2009). The time of divergence between

1111    *Brotula barbata* and Cuban cave brotulas was set to 80 Mya (http://www.timetree.org/).

1112

1113    As an alternative approach, we used the distribution of MutPred2 scores in the lineage leading

1114    to *L. dentata*. First we computed the proportions of two distributions, one under selection as

1115    in the zebrafish lineage ($p_s$) and one without selection as in simulated data ($p_n$), that produce a

1116    mixture distribution that best fit the distribution of MutPred2 scores in the lineage leading to

1117    *L. dentata*. We assumed that $\omega_s$ under selection shifted to $\omega_n$ when selection is relaxed. We

1118    called $T_d$ the period of time since the separation of *L. dentata* and *L gibarensis*, $t_s$ the period of

1119    time of evolution under selection and $t_n$ the period of time under relaxed selection in the

1120    lineage leading to *L. dentata*. In this lineage, the proportion of nonsynonymous substitutions

1121    that accumulate under selection depends on $\omega_s$ and $t_s$ and the proportion of nonsynonymous

45

1122  substitutions that accumulate under relaxed selection depends $\omega_n$ and $t_n$. Thus $\frac{p_s}{p_n} = \frac{t_s \omega_s}{t_n \omega_n}$ or

1123  $t_n = \frac{\omega_s}{\omega_n} \frac{p_n}{p_s} t_s$

1124

1125  ## Comparison of eye gene decay in cavefishes and fossorial mammals

1126

1127  In order to compare the decay of eye genes in cavefishes and fossorial mammals, the number

1128  of pseudogenes and the number of LoF mutations per pseudogene among genes coding for

1129  proteins involved in retinal networks in three fossorial mammals (Cape golden mole

1130  *Chrysochloris asiatica*, naked mole-rat *Heterocephalus glaber* and star-nosed moles

1131  *Condylura cristata*) were retrieved from a publication (Emerling and Springer 2014).

1132

1133  ## Data Availability

1134  *Lucifuga dentata* Whole Genome Shotgun project has been deposited at

1135  DDBJ/ENA/GenBank under the accession VXCM00000000. The version described in this

1136  paper is version VXCM01000000

1137  *Lucifuga dentata* Transcriptome Shotgun Assembly project has been deposited at

1138  DDBJ/EMBL/GenBank under the accession GIAU00000000. The version described in this

1139  paper is the first version, GIAU01000000.

1140  *Lucifuga gibarensis* raw sequences were submitted to the SRA Bioproject: PRJNA610231

1141  The original GFF3 annotation file of *Lucifuga dentata* and scaffolds smaller than 200 bp are

1142  available in Supplementary files.

1143  Python programs and R scripts used in this paper can be found in:

1144  https://github.com/MaximePolicarpo/Molecular-decay-of-light-processing-genes-in-

1145  cavefishes.

1146

## Supplementary Material

1148

1149  Supplementary data are available at Molecular Biology and Evolution.

1150

## Acknowledgments

1152

1159

## Ethics approval

1161

1162  Animals were treated according to the French and European regulations for handling of

1163  animals in research.

1164

## Sampling authorization

1166

1167  *Lucifuga dentata*: a permit [LH 112 AN (135) 2013] was provided to the Centro de

1168  Investigaciones Marinas, University of Havana by the Cuban authorities in December 2013 to

1169  study cave species diversity including nematodes, crustaceans and fishes. As the species was

47

1170    listed Vulnerable (VU) by the IUCN, only two adult individuals (MFP 18.000278) were

1171    sampled (12 January 2014) from one of its largest and demographically stable populations

1172    (Emilio Cave, Las Cañas, Artemisa Province, Cuba).

1173    *Lucifuga gibarensis*: a permit [PE 2014/82] was provided to the Centro de Investigaciones

1174    Marinas, University of Havana by the Cuban authorities in November 2014 to study cave

1175    species diversity including nematodes, crustaceans and fishes. A single adult fish (MFP

1176    18.000279) was sampled (20 November 2014) from the Macigo Cave (Aguada de Macigo del

1177    Jobal), Gibara, Holguín Province, Cuba.

1178

1179

# References

Albalat R, Cañestro C. 2016. Evolution by gene loss. Nature Reviews Genetics 17:379-391.

Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. J Mol Biol 215:403-410.

Audo I, Kohl S, Leroy BP, Munier FL, Guillonneau X, Mohand-Saïd S, Bujakowska K, Nandrot EF, Lorenz B, Preising M, et al. 2009. TRPM1 is mutated in patients with autosomal-recessive complete congenital stationary night blindness. American Journal of Human Genetics 85:720-729.

Beale A, Guibal C, Tamai TK, Klotz L, Cowen S, Peyric E, Reynoso VH, Yamamoto Y, Whitmore D. 2013. Circadian rhythms in Mexican blind cavefish *Astyanax mexicanus* in the lab and in the field. Nature Communications 4:2769.

Benson G. 1999. Tandem repeats finder: a program to analyze DNA sequences. Nucleic Acids Research 27:573-580.

Berning D, Adams H, Luc H, Gross JB. 2019. In-Frame Indel Mutations in the Genome of the Blind Mexican Cavefish, *Astyanax mexicanus*. Genome Biol Evol 11:2563-2573.

Betancur P, Bronner-Fraser M, Sauka-Spengler T. 2010. Assembling Neural Crest Regulatory Circuits into a Gene Regulatory Network. Annual Review of Cell and Developmental Biology 26:581-603.

Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W. 2011. Scaffolding pre-assembled contigs using SSPACE. Bioinformatics 27:578-579.

Burgoyne T, Connor MN, Seabra MC, Cutler DF, Futter CE. 2015. Regulation of melanosome number, shape and movement in the zebrafish retinal pigment epithelium by OA1 and PMEL. Journal of Cell Science 128:1400-1407.

Ceinos RM, Frigato E, Pagano C, Fröhlich N, Negrini P, Cavallari N, Vallone D, Fuselli S, Bertolucci C, Foulkes NS. 2018. Mutations in blind cavefish target the light-regulated circadian clock gene, period 2. Scientific Reports 8:8754.

Chang B, Hawes NL, Hurd RE, Davisson MT, Nusinowitz S, Heckenlively JR. 2002. Retinal degeneration mutants in the mouse. Vision Res 42:517-525.

Chen X, Sheng X, Zhuang W, Sun X, Liu G, Shi X, Huang G, Mei Y, Li Y, Pan X, et al. 2017. GUCA1A mutation causes maculopathy in a five-generation family with a wide spectrum of severity. Genetics in Medicine 19:945-954.

Chikhi R, Rizk G. 2013. Space-efficient and exact de Bruijn graph representation based on a Bloom filter. Algorithms for Molecular Biology 8:22.

Culver DC, Pipan T. 2009. The Biology of Caves and Other Subterranean Habitats. Oxford: Oxford University Press.

David L, Blum S, Feldman MW, Lavi U, Hillel J. 2003. Recent Duplication of the Common Carp (*Cyprinus carpio* L.) Genome as Revealed by Analyses of Microsatellite Loci. Molecular Biology and Evolution 20:1425-1434.

Davies WIL, Tamai TK, Zheng L, Fu JK, Rihel J, Foster RG, Whitmore D, Hankins MW. 2015. An extended family of novel vertebrate photopigments is widely expressed and displays a diversity of function. Genome Research 25:1666-1679.

Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras TR. 2013. STAR: ultrafast universal RNA-seq aligner. Bioinformatics 29:15-21.

Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. Nucleic Acids Research 32:1792-1797.

Eilbeck K, Moore B, Holt C, Yandell M. 2009. Quantitative measures for the management and comparison of annotated genomes. BMC Bioinformatics 10:67.

Emerling CA. 2018. Regressed but Not Gone: Patterns of Vision Gene Loss and Retention in Subterranean Mammals. Integr Comp Biol 58:441-451.

1228 Emerling CA, Springer MS. 2014. Eyes underground: Regression of visual protein networks in
1229 subterranean mammals. Molecular Phylogenetics and Evolution 78:260-270.
1230 Fang X, Nevo E, Han L, Levanon EY, Zhao J, Avivi A, Larkin D, Jiang X, Feranchuk S, Zhu Y, et al. 2014.
1231 Genome-wide adaptive complexes to underground stresses in blind mole rats Spalax. Nature
1232 Communications 5:3966.
1233 Fang X, Seim I, Huang Z, Gerashchenko Maxim V, Xiong Z, Turanov Anton A, Zhu Y, Lobanov Alexei V,
1234 Fan D, Yim Sun H, et al. 2014. Adaptations to a Subterranean Environment and Longevity Revealed by
1235 the Analysis of Mole Rat Genomes. Cell Reports 8:1354-1364.
1236 Farber DB, Lolley RN. 1974. Cyclic Guanosine Monophosphate: Elevation in Degenerating
1237 Photoreceptor Cells of the C3H Mouse Retina. Science 186:449-451.
1238 Feng D, Chen Z, Yang K, Miao S, Xu B, Kang Y, Xie H, Zhao C. 2017. The cytoplasmic tail of rhodopsin
1239 triggers rapid rod degeneration in kinesin-2 mutants. Journal of Biological Chemistry 292:17375-
1240 17386.
1241 Florea L, Souvorov A, Kalbfleisch TS, Salzberg SL. 2011. Genome Assembly Has a Major Impact on
1242 Gene Content: A Comparison of Annotation in Two *Bos Taurus* Assemblies. PLoS ONE 6:e21400.
1243 Fries R, Scholten A, Säftel W, Koch K-W. 2013. Zebrafish Guanylate Cyclase Type 3 Signaling in Cone
1244 Photoreceptors. PLoS ONE 8:e69656.
1245 Fukamachi S, Yada T, Meyer A, Kinoshita M. 2009. Effects of constitutive expression of somatolactin
1246 alpha on skin pigmentation in medaka. Gene 442:81-87.
1247 Fumey J, Hinaux H, Noirot C, Thermes C, Rétaux S, Casane D. 2018. Evidence for late Pleistocene
1248 origin of *Astyanax mexicanus* cavefish. Bmc Evolutionary Biology 18:43.
1249 Gal A, Orth U, Baehr W, Schwinger E, Rosenberg T. 1994. Heterozygous missense mutation in the rod
1250 cGMP phosphodiesterase β–subunit gene in autosomal dominant stationary night blindness. Nature
1251 Genetics 7:64-68.
1252 García-Machado E, Hernandez D, Garcia-Debras A, Chevalier-Monteagudo P, Metcalfe C, Bernatchez
1253 L, Casane D. 2011. Molecular phylogeny and phylogeography of the Cuban cave-fishes of the genus
1254 *Lucifuga*: evidence for cryptic allopatric diversity. Mol Phylogenet Evol 61:470-483.
1255 Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L, Raychowdhury R,
1256 Zeng Q, et al. 2011. Full-length transcriptome assembly from RNA-Seq data without a reference
1257 genome. Nature Biotechnology 29:644.
1258 Grantham R. 1974. Amino acid difference formula to help explain protein evolution. Science 185:862-
1259 864.
1260 Gregory TR. 2019. Animal Genome Size. http://www.genomesize.com.
1261 Gross JB, Borowsky R, Tabin CJ. 2009. A novel role for Mc1r in the parallel evolution of
1262 depigmentation in independent populations of the cavefish *Astyanax mexicanus*. PLoS Genet
1263 5:e1000326.
1264 Gross JB, Weagley J, Stahl BA, Ma L, Espinasa L, McGaugh SE. 2017. A local duplication of the
1265 Melanocortin receptor 1 locus in Astyanax. Genome 61:254-265.
1266 Hinaux H, Blin M, Fumey J, Legendre L, Heuze A, Casane D, Retaux S. 2015. Lens Defects in *Astyanax*
1267 *mexicanus* Cavefish: Evolution of Crystallins and a Role for alphaA-Crystallin. Developmental
1268 Neurobiology 75:505-521.
1269 Hoang DT, Chernomor O, von Haeseler A, Minh BQ, Vinh LS. 2018. UFBoot2: Improving the Ultrafast
1270 Bootstrap Approximation. Molecular Biology and Evolution 35:518-522.
1271 Hoff KJ, Lange S, Lomsadze A, Borodovsky M, Stanke M. 2016. BRAKER1: Unsupervised RNA-Seq-
1272 Based Genome Annotation with GeneMark-ET and AUGUSTUS. Bioinformatics 32:767-769.
1273 Holt C, Yandell M. 2011. MAKER2: an annotation pipeline and genome-database management tool
1274 for second-generation genome projects. BMC Bioinformatics 12:491.
1275 Hunt M, Kikuchi T, Sanders M, Newbold C, Berriman M, Otto TD. 2013. REAPR: a universal tool for
1276 genome assembly evaluation. Genome Biology 14:R47.
1277 Imanishi Y, Li N, Sokal I, Sowa ME, Lichtarge O, Wensel TG, Saperstein DA, Baehr W, Palczewski K.
1278 2002. Characterization of retinal guanylate cyclase-activating protein 3 (GCAP3) from zebrafish to
1279 man. European Journal of Neuroscience 15:63-78.

1280    Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell A, Nuka G,
1281    et al. 2014. InterProScan 5: genome-scale protein function classification. Bioinformatics 30:1236-
1282    1240.
1283    Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermiin LS. 2017. ModelFinder: fast model
1284    selection for accurate phylogenetic estimates. Nature Methods 14:587-589.
1285    Kastenhuber E, Gesemann M, Mickoleit M, Neuhauss SCF. 2013. Phylogenetic analysis and expression
1286    of zebrafish transient receptor potential melastatin family genes. Developmental Dynamics
1287    242:1236-1249.
1288    Kim EB, Fang X, Fushan AA, Huang Z, Lobanov AV, Han L, Marino SM, Sun X, Turanov AA, Yang P, et al.
1289    2011. Genome sequencing reveals insights into physiology and longevity of the naked mole rat.
1290    Nature 479:223-227.
1291    Kimura T, Nagao Y, Hashimoto H, Yamamoto-Shiraishi Y-i, Yamamoto S, Yabe T, Takada S, Kinoshita
1292    M, Kuroiwa A, Naruse K. 2014. Leucophores are similar to xanthophores in their specification and
1293    differentiation processes in medaka. Proceedings of the National Academy of Sciences of the United
1294    States of America 111:7343-7348.
1295    Kohany O, Gentles AJ, Hankus L, Jurka J. 2006. Annotation, submission and screening of repetitive
1296    elements in Repbase: RepbaseSubmitter and Censor. BMC Bioinformatics 7:474.
1297    Krauss J, Frohnhöfer HG, Walderich B, Maischein H-M, Weiler C, Irion U, Nüsslein-Volhard C. 2014.
1298    Endothelin signalling in iridophore development and stripe pattern formation of zebrafish. Biology
1299    Open 3:503-509.
1300    Krauss J, Geiger-Rudolph S, Koch I, Nüsslein-Volhard C, Irion U. 2014. A dominant mutation in tyrp1A
1301    leads to melanophore death in zebrafish. Pigment Cell & Melanoma Research 27:827-830.
1302    Kriventseva EV, Zdobnov EM, Simão FA, Ioannidis P, Waterhouse RM. 2015. BUSCO: assessing
1303    genome assembly and annotation completeness with single-copy orthologs. Bioinformatics 31:3210-
1304    3212.
1305    Lagman D, Callado-Pérez A, Franzén IE, Larhammar D, Abalo XM. 2015. Transducin Duplicates in the
1306    Zebrafish Retina and Pineal Complex: Differential Specialisation after the Teleost Tetraploidisation.
1307    PLoS ONE 10:e0121330.
1308    Lagman D, Franzén IE, Eggert J, Larhammar D, Abalo XM. 2016. Evolution and expression of the
1309    phosphodiesterase 6 genes unveils vertebrate novelty to control photosensitivity. Bmc Evolutionary
1310    Biology 16:124.
1311    Lahti DC, Johnson NA, Ajie BC, Otto SP, Hendry AP, Blumstein DT. 2009. Relaxed selection in the wild.
1312    Trends Ecol Evol 24:487–496.
1313    Letunic I, Bork P. 2006. Interactive Tree Of Life (iTOL): an online tool for phylogenetic tree display and
1314    annotation. Bioinformatics 23:127-128.
1315    Li H. 2011. A statistical framework for SNP calling, mutation discovery, association mapping and
1316    population genetical parameter estimation from sequencing data. Bioinformatics 27:2987-2993.
1317    Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
1318    Bioinformatics 25:1754-1760.
1319    Li W-H, Gojobori T, Nei M. 1981. Pseudogenes as a paradigm of neutral evolution. Nature 292:237-
1320    239.
1321    Li W-H, Nei M. 1977. Persistence of Common Alleles in Two Related Populations or Species. Genetics
1322    86:901-914.
1323    Li Y, Li G, Wang H, Du J, Yan J. 2013. Analysis of a Gene Regulatory Cascade Mediating Circadian
1324    Rhythm in Zebrafish. Plos Computational Biology 9:e1002940.
1325    Lin J-J, Wang F-Y, Li W-H, Wang T-Y. 2017. The rises and falls of opsin genes in 59 ray-finned fish
1326    genomes and their implications for environmental adaptation. Scientific Reports 7:15568.
1327    Lister JA. 2019. Larval but not adult xanthophore pigmentation in zebrafish requires GTP
1328    cyclohydrolase 2 (gch2) function. Pigment Cell & Melanoma Research 0.
1329    Liu Z, Wen H, Hailer F, Dong F, Yang Z, Liu T, Han L, Shi F, Hu Y, Zhou J. 2019. Pseudogenization of
1330    *Mc1r* gene associated with transcriptional changes related to melanogenesis explains leucistic

1331    phenotypes in *Oreonectes* cavefish (Cypriniformes, Nemacheilidae). Journal of Zoological Systematics
1332    and Evolutionary Research 57:900-909.
1333    Lorin T, Brunet FG, Laudet V, Volff J-N. 2018. Teleost Fish-Specific Preferential Retention of
1334    Pigmentation Gene-Containing Families After Whole Genome Duplications in Vertebrates. G3:
1335    Genes|Genomes|Genetics 8:1795-1806.
1336    Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Pan Q, Liu Y, et al. 2012. SOAPdenovo2: an
1337    empirically improved memory-efficient short-read de novo assembler. GigaScience 1:18.
1338    Lynch M, Conery JS. 2000. The evolutionary fate and consequences of duplicate genes. Science
1339    290:1151-1155.
1340    Lynch M, Kewalramani A. 2003. Messenger RNA Surveillance and the Evolutionary Proliferation of
1341    Introns. Molecular Biology and Evolution 20:563-571.
1342    Ma L, Parkhurst A, Jeffery W. 2014. The role of a lens survival pathway including sox2 and alphaA-
1343    crystallin in the evolution of cavefish eye degeneration. Evodevo 5:28.
1344    Mackay DS, Boskovska OB, Knopf HLS, Lampi KJ, Shiels A. 2002. A Nonsense Mutation in CRYBB1
1345    Associated with Autosomal Dominant Cataract Linked to Human Chromosome 22q. American Journal
1346    of Human Genetics 71:1216-1221.
1347    Makino CL, Wen X-H, Olshevskaya EV, Peshenko IV, Savchenko AB, Dizhoor AM. 2012. Enzymatic
1348    Relay Mechanism Stimulates Cyclic GMP Synthesis in Rod Photoresponse: Biochemical and
1349    Physiological Study in Guanylyl Cyclase Activating Protein 1 Knockout Mice. PLoS ONE 7:e47637.
1350    Malmstrøm M, Matschiner M, Tørresen OK, Jakobsen KS, Jentoft S. 2017. Whole genome sequencing
1351    data and de novo draft assemblies for 66 teleost species. Scientific data 4:160132.
1352    Mao H, Wang H. 2017. SINE_scan: an efficient tool to discover short interspersed nuclear elements
1353    (SINEs) in large-scale genomic datasets. Bioinformatics 33:743-745.
1354    Matveev AV, Quiambao AB, Fitzgerald JB, Ding XQ. 2008. Native cone photoreceptor cyclic
1355    nucleotide-gated channel is a heterotetrameric complex comprising both CNGA3 and CNGB3: a study
1356    using the cone-dominant retina of Nrl–/– mice. Journal of Neurochemistry 106:2042-2055.
1357    McGaugh SE, Gross JB, Aken B, Blin M, Borowsky R, Chalopin D, Hinaux H, Jeffery WR, Keene A, Ma L,
1358    et al. 2014. The cavefish genome reveals candidate genes for eye loss. Nat Commun 5:5307.
1359    McLaughlin ME, Sandberg MA, Berson EL, Dryja TP. 1993. Recessive mutations in the gene encoding
1360    the β−subunit of rod phosphodiesterase in patients with *retinitis pigmentosa*. Nature Genetics 4:130-
1361    134.
1362    Meredith RW, Gatesy J, Murphy WJ, Ryder OA, Springer MS. 2009. Molecular Decay of the Tooth
1363    Gene Enamelin (ENAM) Mirrors the Loss of Enamel in the Fossil Record of Placental Mammals. Plos
1364    Genetics 5:e1000634.
1365    Miyata T, Yasunaga T. 1981. Rapidly evolving mouse alpha-globin-related pseudo gene and its
1366    evolutionary history. Proceedings of the National Academy of Sciences 78:450-453.
1367    Morgulis A, Gertz EM, Schäffer AA, Agarwala R. 2006. A Fast and Symmetric DUST Implementation to
1368    Mask Low-Complexity DNA Sequences. Journal of Computational Biology 13:1028-1040.
1369    Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: A Fast and Effective Stochastic
1370    Algorithm for Estimating Maximum-Likelihood Phylogenies. Molecular Biology and Evolution 32:268-
1371    274.
1372    Nishiwaki Y, Komori A, Sagara H, Suzuki E, Manabe T, Hosoya T, Nojima Y, Wada H, Tanaka H,
1373    Okamoto H, et al. 2008. Mutation of cGMP phosphodiesterase 6α⍰-subunit gene causes progressive
1374    degeneration of cone photoreceptors in zebrafish. Mechanisms of Development 125:932-946.
1375    Nord H, Dennhag N, Muck J, von Hofsten J. 2016. Pax7 is required for establishment of the
1376    xanthophore lineage in zebrafish embryos. Molecular biology of the cell 27:1853-1862.
1377    Osawa S, Weiss ER editors. Retinal Degenerative Diseases. 2012 Boston, MA.
1378    Parichy DM, Ransom DG, Paw B, Zon LI, Johnson SL. 2000. An orthologue of the kit-related gene fms
1379    is required for development of neural crest-derived xanthophores and a subpopulation of adult
1380    melanocytes in the zebrafish, Danio rerio. Development 127:3031.
1381    Parichy DM, Spiewak JE. 2015. Origins of adult pigmentation: diversity in pigment stem cell lineages
1382    and implications for pattern evolution. Pigment Cell & Melanoma Research 28:31-50.

1383 Pejaver V, Urresti J, Lugo-Martinez J, Pagel KA, Lin GN, Nam H-J, Mort M, Cooper DN, Sebat J,
1384 Iakoucheva LM, et al. 2017. MutPred2: inferring the molecular and phenotypic impact of amino acid
1385 variants. bioRxiv:134981.
1386 Pelletier I, Bally-Cuif L, Ziegler I. 2001. Cloning and developmental expression of zebrafish GTP
1387 cyclohydrolase I. Mechanisms of Development 109:99-103.
1388 Protas ME, Hersey C, Kochanek D, Zhou Y, Wilkens H, Jeffery WR, Zon LI, Borowsky R, Tabin CJ. 2006.
1389 Genetic analysis of cavefish reveals molecular convergence in the evolution of albinism. Nat Genet
1390 38:107-111.
1391 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features.
1392 Bioinformatics 26:841-842.
1393 Rätscho N, Scholten A, Koch K-W. 2009. Expression profiles of three novel sensory guanylate cyclases
1394 and guanylate cyclase-activating proteins in the zebrafish retina. Biochimica et Biophysica Acta (BBA)
1395 - Molecular Cell Research 1793:1110-1114.
1396 Renninger SL, Gesemann M, Neuhauss SCF. 2011. Cone arrestin confers cone vision of high temporal
1397 resolution in zebrafish larvae. European Journal of Neuroscience 33:658-667.
1398 Sato M, Nakazawa M, Usui T, Tanimoto N, Abe H, Ohguro H. 2005. Mutations in the gene coding for
1399 guanylate cyclase-activating protein 2 (GUCA1B gene) in patients with autosomal dominant retinal
1400 dystrophies. Graefe's Archive for Clinical and Experimental Ophthalmology 243:235-242.
1401 Schonthaler HB, Lampert JM, Isken A, Rinner O, Mader A, Gesemann M, Oberhauser V, Golczak M,
1402 Biehlmaier O, Palczewski K, et al. 2007. Evidence for RPE65-independent vision in the cone-
1403 dominated zebrafish retina. The European journal of neuroscience 26:1940-1949.
1404 Simon N, Fujita S, Porter M, Yoshizawa M. 2019. Expression of extraocular opsin genes and light-
1405 dependent basal activity of blind cavefish. PeerJ 7:e8148.
1406 Slater GSC, Birney E. 2005. Automated generation of heuristics for biological sequence comparison.
1407 BMC Bioinformatics 6:31.
1408 Smit AFA, Hubley R. 2008. RepeatModeler Open 1.0 . http://www.repeatmasker.org.
1409 Smit AFA, Hubley R, Green P. 2013. RepeatMasker Open 4.0. http://www.repeatmasker.org.
1410 Stearns G, Evangelista M, Fadool JM, Brockerhoff SE. 2007. A mutation in the cone-specific *pde6* gene
1411 causes rapid cone photoreceptor degeneration in zebrafish. The Journal of Neuroscience 27:13866 –
1412 13874.
1413 Takahashi JS, Hong H-K, Ko CH, McDearmon EL. 2008. The genetics of mammalian circadian order
1414 and disorder: implications for physiology and disease. Nature Reviews Genetics 9:764-775.
1415 Thanos S, Böhm MRR, Meyer zu Hörste M, Prokosch-Willing V, Hennig M, Bauer D, Heiligenhaus A.
1416 2014. Role of crystallins in ocular neuroprotection and axonal regeneration. Progress in Retinal and
1417 Eye Research 42:145-161.
1418 Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL, Wold BJ, Pachter
1419 L. 2010. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and
1420 isoform switching during cell differentiation. Nature Biotechnology 28:511-515.
1421 Vatine G, Vallone D, Gothilf Y, Foulkes NS. 2011. It's time to swim! Zebrafish and the circadian clock.
1422 FEBS Letters 585:1485-1494.
1423 Vincent A, Audo I, Tavares E, Maynes Jason T, Tumber A, Wright T, Li S, Michiels C, Banin E, Bocquet
1424 B, et al. 2016. Biallelic Mutations in *GNB3* Cause a Unique Form of Autosomal-Recessive Congenital
1425 Stationary Night Blindness. The American Journal of Human Genetics 98:1011-1019.
1426 Vurture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H, Gurtowski J, Schatz MC. 2017.
1427 GenomeScope: fast reference-free genome profiling from short reads. Bioinformatics 33:2202-2204.
1428 Wada Y, Sugiyama J, Okano T, Fukada Y. 2006. GRK1 and GRK7: Unique cellular distribution and
1429 widely different activities of opsin phosphorylation in the zebrafish rods and cones. Journal of
1430 Neurochemistry 98:824-837.
1431 Wasfy MM, Matsui JI, Miller J, Dowling JE, Perkins BD. 2014. myosin 7aa−/− mutant zebrafish show
1432 mild photoreceptor degeneration and reduced electroretinographic responses. Experimental Eye
1433 Research 122:65-76.

1434    Wertheim JO, Murrell B, Smith MD, Kosakovsky Pond SL, Scheffler K. 2015. RELAX: Detecting Relaxed
1435    Selection in a Phylogenetic Framework. Molecular Biology and Evolution 32:820-832.
1436    Xiao H, Chen S-y, Liu Z-m, Zhang R-d, Li W-x, Zan R-g, Zhang Y-p. 2005. Molecular phylogeny of
1437    *Sinocyclocheilus* (Cypriniformes: Cyprinidae) inferred from mitochondrial DNA sequences. Molecular
1438    Phylogenetics and Evolution 36:67-77.
1439    Yang J, Chen X, Bai J, Fang D, Qiu Y, Jiang W, Yuan H, Bian C, Lu J, He S, et al. 2016. The
1440    *Sinocyclocheilus* cavefish genome provides insights into cave adaptation. BMC Biol 14:1-13.
1441    Yang Z. 2007. PAML 4: Phylogenetic Analysis by Maximum Likelihood. Molecular Biology and
1442    Evolution 24:1586-1591.
1443    Yuan J, He Z, Yuan X, Jiang X, Sun X, Zou S. 2010. Speciation of polyploid Cyprinidae fish of common
1444    carp, crucian carp, and silver crucian carp derived from duplicated Hox genes. Journal of
1445    Experimental Zoology Part B: Molecular and Developmental Evolution 314B:445-456.
1446    Zang J, Keim J, Kastenhuber E, Gesemann M, Neuhauss SCF. 2015. Recoverin depletion accelerates
1447    cone photoresponse recovery. Open Biology 5.
1448    Zhao H, Di Mauro G, Lungu-Mitea S, Negrini P, Guarino AM, Frigato E, Braunbeck T, Ma H, Lamparter
1449    T, Vallone D, et al. 2018. Modulation of DNA Repair Systems in Blind Cavefish during Evolution in
1450    Constant Darkness. Current Biology 28:3229-3243.
1451    Zhao H, Yang J-R, Xu H, Zhang J. 2010. Pseudogenization of the Umami Taste Receptor Gene Tas1r1 in
1452    the Giant Panda Coincided with its Dietary Switch to Bamboo. Molecular Biology and Evolution
1453    27:2669-2673.

1454

1455

1456    Legends

1457

1458    **Fig. 1.** Gene sets. (A) Eye genes. (B) Circadian clock genes. (C) Pigmentation genes. Genes

1459    were colored according to the species in which LoF mutations were found (species name

1460    followed by * indicate that no genome was available but pseudogenes were identified). Genes

1461    with a hatched background are under/not expressed in at least one cavefish species. Genes

1462    expressed only in the eyes are surrounded by blue lines and opsins by an orange line.

1463    Pigmentation genes were clustered according to Lorin et al. 2018. Genes belonging to several

1464    subsets are surrounded by dotted lines with links between the different subsets.

1465

1466    **Fig. 2.** Phylogeny of the cavefishes and some close relatives.

1467

1468    **Fig. 3.** Mapping of LoF mutations. For *Sinocyclocheilus* species, the number of genes for

1469    which both ohnologs are pseudogenized is given between brackets.

1470

1471    **Fig. 4.** Distribution of different categories of LoF mutations. (A) Position of internal stop

1472    codons and frameshifts along coding sequences. (B) Observed and expected frequencies. (C)

1473    Distribution of indel size.

1474

1475    **Fig. 5.** Observed and expected distributions of LoF mutations per gene. (A) *L. dentate*. (B) *L.*

1476    *gibarensis*. Red line: observed distribution. The expected distributions were obtained using an

1477    analytical model (dots) and 10,000 simulations (histograms).

1478

1479    **Fig. 6.** Distribution of ω in surface and cave fishes. (A) Eye genes. (B) Circadian clock genes.

1480    (C) Pigmentation genes.

1481

1482 **Fig. 7.** Distributions of MutPred2 scores in several fish lineages and in simulations of

1483 substitutions without selection. The number of substitutions in each lineage is given between

1484 parenthesis. One hundred simulations were performed with each gene set. In each simulation

1485 54 non-synonymous mutations were generated in eye genes, 36 in circadian clock genes and

1486 232 in pigmentation genes, those numbers corresponding to the numbers of non-synonymous

1487 mutations found in *Astyanax mexicanus* cavefish.

1488

1489 **Fig. 8.** Probability of finding 19 eye pseudogenes in *L. dentata* according to the time of

1490 neutral evolution. Red and pink lines: based on an analytical model assuming 76 and 50

1491 neutral genes respectively; other lines: estimations based on 10,000 simulations, assuming 76

1492 neutral genes and taking into account the length and number of introns in each eye genes and

1493 considering different effective population sizes. The number of generations for which the

1494 highest probability was found is reported above each line.

1495

1496 **Fig. 9.** Comparison of eye gene decay in cavefishes *vs* fossorial mammals.

1497 (A) Venn diagram showing the genes carrying LoF mutations in both groups. For each gene,

1498 the number of LoF mutations found in each species is indicated. (B) Distribution of the

1499 number of LoF mutations per pseudogene. The distribution was computed with only the

1500 pseudogenes found in both groups or with all pseudogenes (inset). Genes present as one copy

1501 in fossorial mammals are often duplicated in *L. dentata* and quadruplicated in *S. anshuiensis,*

1502 after one and two WGD respectively. Other gene duplications also sporadically increased the

1503 number of paralogs in these fishes. The number of LoF mutations found in these paralogs are

1504 separated by vertical lines.

1505

1506 **Table1.** Estimations of the time without selection on eye genes in *Lucifuga dentata*.

1507

1508 **Fig. S1.** Photos of the specimens used for genome sequencing. (A) *Lucifuga dentate*. (B)

1509 *Lucifuga gibarensis*.

1510

1511 **Fig. S2.** *L. dentata* scaffold size distribution. Based on 3,537 scaffolds longer than 50kb

1512 (49,407 scaffolds < 50 kb not used).

1513

1514 **Fig. S3.** BUSCO analyses using the Actinopterygii gene database (v3.1.0). We assessed the

1515 completeness of three published Ophidiiformes genomes (*Brotula barbata*, *Carapus acus* and

1516 *Lamprogrammus exutus*), *Lucifuga dentata* genome, gene models resulting from the

1517 annotation pipeline and transcriptome assembly.

1518

1519 **Fig. S4.** Transcriptome statistics. We followed the transcriptome assembly quality assessment

1520 of Trinity (https://github.com/trinityrnaseq/trinityrnaseq/wiki/Transcriptome-Assembly-

1521 Quality-Assessment).

1522

1523 **Fig. S5.** Genome annotation pipeline used on the *Lucifuga dentata* draft genome.

1524

1525 **Fig. S6.** Interspersed repeat landscape of *Lucifuga dentata*.

1526

1527 **Fig. S7.** List of eye genes retrieved from cavefishes and related species. Colors represent the

1528 type of LoF mutation. When higher than one, the number of LoF mutations is also reported.

1529

1530 **Fig. S8.** A large deletion between GNL3L and SWS2 at the origin of LWS gene loss in

1531 Ophidiiformes. This figure was generated using SimpleSynteny (Veltri D., Malapi-Wight M.

1532 and Crouch J.A. SimpleSynteny: a web-based tool for visualization of microsynteny across

1533 multiple species. *Nucleic Acids Research* 44(W1):W41-W45, 2016, doi:10.1093/nar/gkw330).

1534

1535 **Fig. S9.** Distribution of the effective segment size generated by random insertion of STOP

1536 codons and frameshifts (100,000 simulations).

1537

1538 **Fig. S10.** Frameshift size distribution for each dataset.

1539

1540 **Fig. S11.** Observed and theoretical frequencies of different types of LoF mutations in three

1541 gene sets, and the frequency of different types of mutations found in eye genes of fossorial

1542 mammals (Emerling CA, Springer MS. 2014. Eyes underground: Regression of visual protein

1543 networks in subterranean mammals. Molecular Phylogenetics and Evolution 78:260-270).

1544

1545 **Fig. S12.** (A) Number of difference per gene between *Astyanax mexicanus* morphs and

1546 between *Lucifuga dentata* and *Lucifuga gibarensis*. (B) Estimation of ω for each eye gene.

1547 Grey lines represent values of dn or ds < 0.01, leading to non-reliable estimations of ω.

1548

1549 **Fig. S13**. Estimations of ω with concatenated sequences. Branch colors are scaled depending

1550 on the ω values. Trees were generated using ggtree (Yu, G., Smith, D.K., Zhu, H., Guan, Y.

1551 and Lam, T.T.-Y. (2017), ggtree: an R package for visualization and annotation of

1552 phylogenetic trees with their covariates and other associated data. Methods Ecol Evol, 8: 28-

1553 36. doi:10.1111/2041-210X.12628).

1554

1555    **Fig. S14.** RELAX results with species assigned as test branch for eye genes. The k parameter

1556    and p-value are displayed along with ω plots.

1557

1558    **Fig. S15.** RELAX results with species assigned as test branch for circadian clock genes. The k

1559    parameter and p-value are displayed along with ω plots.

1560

1561    **Fig. S16.** RELAX results with species assigned as test branch for pigmentation genes. The k

1562    parameter and p-value are displayed along with ω plots.

1563

1564    **Fig. S17.** Empirical cumulative distributions of MutPred2 scores. The number of scores is

1565    indicated between parenthesis. 100 neutral simulations were performed for each dataset with

1566    54 random non synonymous mutations in eye genes, 36 in circadian clock genes and 232 in

1567    pigmentation genes, which are the number of non-synonymous mutations found in *Astyanax*

1568    *mexicanus* cavefish. The statistical significance of the difference between each pair of

1569    distributions was assessed using the Kolmogorov-Smirnov test (significant differences are

1570    shown on a red background whereas non-significant differences are shown on a green

1571    background).

1572

1573    **Fig. S18.** Empirical cumulative distributions of Grantham's distances. The number of

1574    distances is indicated between parenthesis. 100 neutral simulations were performed for each

1575    dataset with 54 random non synonymous mutations in eye genes, 36 in circadian clock genes

1576    and 232 in pigmentation genes which are the number of non-synonymous mutations found in

1577    *Astyanax mexicanus* cavefish. The statistical significance of the difference between each pair

1578    of distributions was assessed using the Kolmogorov-Smirnov test (significant differences are

1579     shown on a red background whereas non-significant differences are shown on a green

1580     background).

1581

1582     **Fig. S19.** Effect of the Transition/Transversion ratio on the cumulative distribution of

1583     MutPred2 scores in simulated amino acid substitutions.

1584

1585     **Fig. S20.** Fit of mixture distributions of MutPred2 scores with the distributions found in two

1586     *Lucifuga* spp. and two *Astyanax mexicanus* morphs. The p-values of Kolomogorv-Smirnov

1587     tests between the observed distributions in each species and mixture distributions were plotted

1588     according to different proportions of mutations that reached fixation under relaxed selection.

1589

1590     **Fig. S21.** Distributions of MutPred2 scores in three *Sinocyclocheilus* species, *Danio rerio*

1591     and in simulations of substitutions without selection. The number of substitutions in each

1592     lineage is given between parenthesis. One hundred simulations were performed with each

1593     gene set. In each simulation 54 non-synonymous mutations were generated in eye genes, 36 in

1594     circadian clock genes and 232 in pigmentation genes, those numbers corresponding to the

1595     numbers of non-synonymous mutations found in *Astyanax mexicanus* cavefish.

1596

1597     **Fig. S22.** Fit of mixture distributions of MutPred2 scores with the distributions found in three

1598     *Sinocyclocheilus* species. The p-values of Kolomogorv-Smirnov tests between the observed

1599     distributions in each species and mixture distributions were plotted according to different

1600     proportions of mutations that reached fixation under relaxed selection.

1601

1602     **Table S1.** List of LoF mutations found in *Lucifuga dentata* and *Lucifuga gibarensis* genomes,

1603     and their coverage. LoF mutations in red were also found in the transcriptome of *L. dentata*.

60

1604

1605    **Data_Supp1.** Summary of the number of genes retrieved from each species and for each gene

1606    set, along with the number of pseudogenes and the number of LoF mutations.

1607

1608    **Data_Supp2.** Sequences predicted with exonerate and ID of sequences retrieved from

1609    Ensembl.

1610

1611    **Data_Supp3.** Results obtained with different methods for dating relaxed selection on eye

1612    genes in *Lucifuga dentata*.

1613

1614    Description of Supplementary files content:

1615    Divergence_values: Pairwise nucleotidic distances between species for each gene set.

1616    Lucifuga_Supplementary_files_Genome: Original GFF3 file with functional annotations and

1617    scaffolds smaller than 200 bp not uploaded to NCBI.

1618    MutPred2_Results: Raw output of MutPred2. Parsed results files to be used with the script

1619    provided in github (MutPred2_Script.R) are also provided.

1620    Phylogenies: Gene phylogenies computed with iQTree and displayed with iTOL. The model

1621    used for each phylogeny can be found on the "Models" folder.

1622    Concatenated_Alignments: Concatenated alignments for vision, circadian and pigmentation

1623    genes.

A

**Transducins**

gnat2 gna13a gnb3b gngt2a
gnat1 gnb1a gnb1b
gngt1 gngt2b

**PDE6 complex**

pde6a pde6b pde6c
pde6ga pde6gb pde6ha pde6hb

**Visual opsins**

rh1.1 rh1.2

rh2.1 rh2.2 rh2.3 rh2.4
lws1 lws2 sws1 sws2

**Non-visual opsins**

opn4m1 opn4m2 opn4m3 opn4x1 opn4x2 opn5 opn6a
opn6b opn7a opn7b opn7c opn7d opn8a opn8b
opn8c opn9 Para-1 Para-2 Pariet. rgr1 rgr2
rrh tmt1a tmt1b tmt2a tmt2b tmt3a tmt3b
exorh opn3 va1 va2

**Guanylate cyclases**

gc2 gc3 gucy2f

**GC-activating proteins**

gcap3 gcap4 gcap5 gcap7
gcap1 gcap2

**Arrestins**

saga arr3a arr3b sagb

**Beta-crystallins**

cryba1b cryba1l1
cryba2a cryba2b
cryba4 crybb1
crybb1l1 crybb1l2
crybb1l3 crybgx

**Alpha-crystallin**

cryaa

**Gamma-crystallins**

crygm5 crygn2

**Recoverins**

rcv1a rcv1b rcv2a rcv2b

**RPE**

rpe65a
rpe65b/c

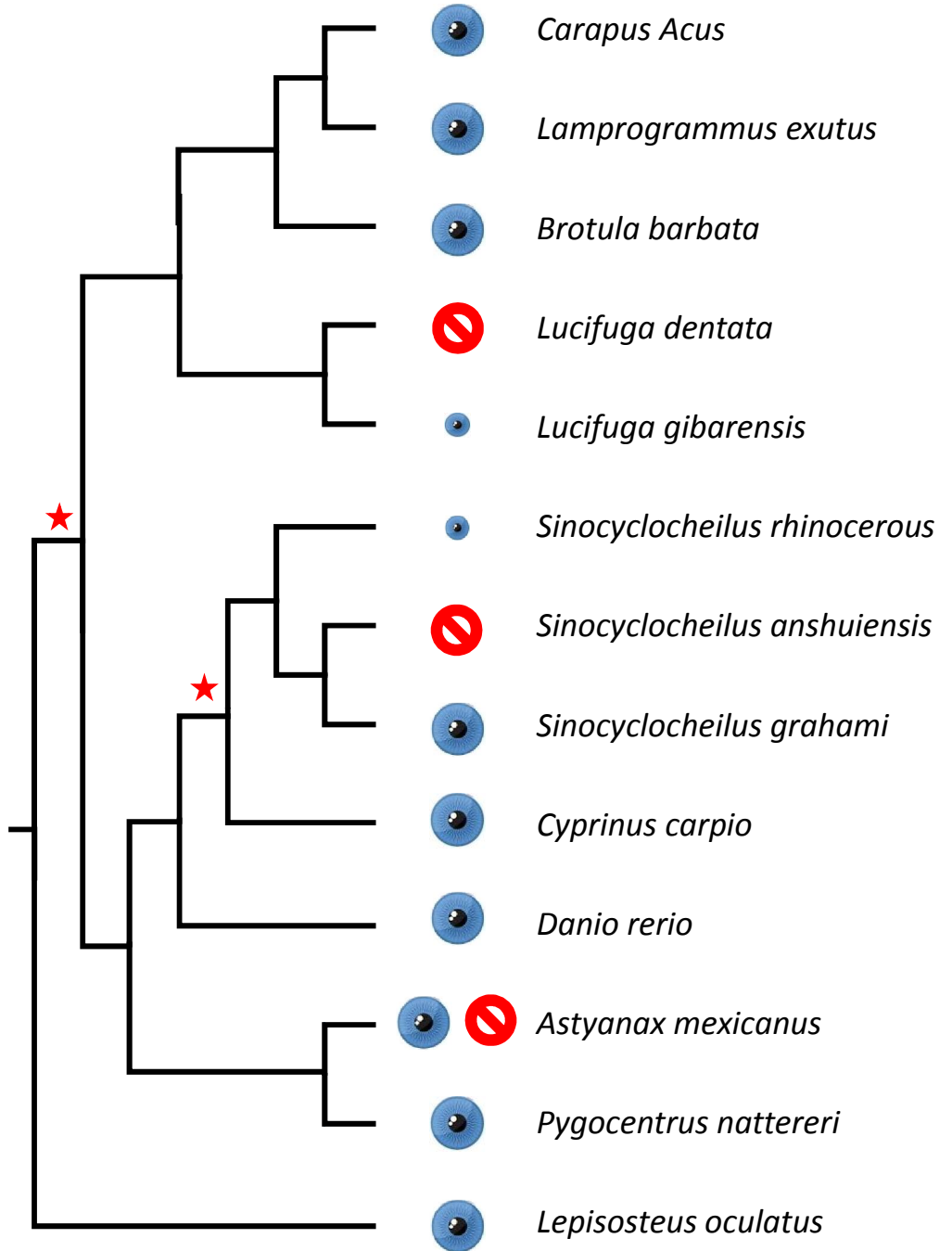**G-protein-coupled receptor kinases**

grk1a grk1b grk7a grk7b

Cone gene    Cone + Rod gene    Rod gene

B

**Brain and Muscle ARNT-Like**

bmal1a bmal1b bmal2

**Clock circadian regulator**

clock1 clock2 clock3

**Cryptochrome**

cry1a cry1b cry2a cry2b
cry3 cry4 cry5 dash

**Period**

per1a per1b
per2 per3

**Casein kinase**

ck1δa ck1δb

**D site albumin promoter binding protein**

dbpa dbpb

**HLF transcription factor**

hlfa hlfb

**TEF transcription factor**

tefa tefb

**Nuclear factor, interleukin 3 regulated**

nfil3 nfil3-2 nfil3-4 nfil3-6

**Nuclear receptor, group D**

nr1d1 nr1d2a nr1d2b nr1d4a nr1d4b

**arylalkylamine N-acetyltransferase**

aanat1 aanat2

**RAR-related orphan receptor**

roraa rorab rorb rorca rorcb

*Lucifuga dentata*    *Lucifuga gibarensis*    *Astyanax mexicanus* (Pachón)

*Sinocyclocheilus anshuiensis*    *Sinocyclocheilus grahami*    *Sinocyclocheilus rhinocerous*

*Amblyopsidae* cavefishes spp.*    *Phreatichthys andruzzii* *    *Oreonectes daqikongensis* *

C

**Pigment cell differentiation**

cdh2  lef1
nrarpa  nrarpb
ovol1a  ovol1b

bnc2
wnt1
wnt3a

ednrba  ednrbb
oca2  tfap2a
sox9a  sox9b
impdh1a
impdh1b

**Iridrophores**

fhl2a  fhl2b  chm
pnp4a  csf1b  ece2a  pnp4b
csf1a  ece2b  mpv17  foxd3
trim33  med12  fbxw4  ltk

**Melanophore development**

adam17a  adamts20  adam17b
bcl2a  bcl2b  adrb2a  adrb2b  atp6v0b
apc  bcl2l11  psen1  fgfr2  brsk2a  brsk2b  cited1
c10orf11  creb1a  creb1b  fzd4  gli3  dct  dock7  ece1
eda  en1a  en1b  edar  edn3a  edn3b  ednrb2  gnaq  gfpt1
frem2a  frem2b  egfr  erbb3a  erbb3b  gata3  gja5a  gja5b  gpc3  gpr161
gna11a  gna11b  sox2  hdac1  hps1  hps4  hps6  irf4a  irf4b  ikbkg
igsf11  hsd3b1  kitlga  kitlgb  lmx1a  mbtps1  sox18  tfap2e  nrg1  oprm1
kcnj13  kita  kitb  mef2ca  mef2cb  mib1  mib2  mcoln3a  mcoln3b
mitfa  mitfb  mreg  myca  mycb  myo5aa  myo5ab  rnf41  rab27a
itgb1a  itgb1b  itgb1b.1  itgb1b.2  scarb2a  scarb2b  scg2a  scg2b
recql4  sf3b1  sfxn1  mtrex  slc24a5  slc45a2  usp13
snai2  trpm7  tyrp1a  tyrp1b  zic2a  zic2b
vps11  vps18  tyr

sox10

gch2
csf1ra
csf1rb

**Xantophores and leucophores development**

ghrb  leo1
gch1  xdh  ghra  mycbp2  sox5
slc2a15a  slc2a15b  smtla  smtlb  gchfr
slc2a11a  slc2a11b  pcbd1  pcbd2  pax3a  pts
slc2a11l  pax7a  qdprb1  qdpra  pax3b
pax7b  qdprb2  spra  sprb

bloc1s6

**Melanosome transport**

ippk  dctn2  tmem33  myo7aa  map2k1
rab1aa  rab17  rab1ab  ric8b  agrp2
crha  crhb  dctn1a  dctn1b  rab11a
mlpha  rab3ip  mlphb  rab8a
rab11ba  rab11bb  myo7ab

tpcn2

**Melanogenesis regulation**

asip1  atrn  myg1  clcn7  corin  ctns
drd2a  drd2b  mc1r  krt2a  defb  shroom2a
slc7a11  ugt1a8  nf1a  nf1b  ostm1  shroom2b
mgrn1a  pah  pomca  pomcb  zeb2a  zeb2b

gart
paics

pmela
pmelb
trappc6a

rab32a  rab32b
rab38a  rab38b

**Melanosome biogenesis**

ankrd27  ap1g1  ap1m1  ap3b1  ap3d1  arcn1a  arcn1b
bloc1s1  bloc1s2  bloc1s3  bloc1s4  bloc1s5  rabggta  fig4
dtnbp1a  dtnbp1b  hps3  hps5  cd63  gpr143  th
kif13a  lyst  mlana  vps33a  vps39  txndc5  snapin
nsfa  nsfb

**Components of melanosomes**

gpnmb  vat1  sdcbp  trpm1a
slc24a4a  slc24a4b  tspan10  trpm1b

**A**

rgr1

adamts20

$\frac{19}{76}$ $\frac{0}{36}$ $\frac{8}{237}$

*Lucifuga dentata*

$\frac{5}{76}$ $\frac{0}{36}$ $\frac{7}{237}$

*Lucifuga gibarensis*

**B**

sws1

gc3/sws1

grk7b
crygm5

pmelb

ovol1b
recql4
mlana

cry-d

$\frac{48[7]}{171}$ $\frac{9[1]}{81}$ $\frac{35[2]}{487}$

*Sinocyclocheilus anshuiensis*

$\frac{18[1]}{173}$ $\frac{5[0]}{80}$ $\frac{15[0]}{484}$

*Sinocyclocheilus grahami*

$\frac{32[1]}{169}$ $\frac{15[4]}{83}$ $\frac{28[2]}{490}$

*Sinocyclocheilus rhinocerous*

*Cyprinus carpio*

**C**

$\frac{1}{85}$ $\frac{0}{38}$ $\frac{2}{249}$

*Astyanax mexicanus* (Pachón)

$\frac{0}{86}$ $\frac{0}{38}$ $\frac{0}{249}$

*Astyanax mexicanus* (Surface)

🛑 Gain of stop codon       ▲ Insertion       🟢 Loss of start codon       | Vision |       Number of pseudogenes / Number of genes

◆ Splice site mutation       ▼ Deletion       🟠 Loss of stop codon       | Pigmentation | | Circadian clock |

# A

Position of stop codons on CDS (length in %)

0                                                                    100

Position of frameshifts on CDS (length in %)

0                                                                    100

● Eye    ● Circadian clock    ● Pigmentation

# B

**Expected LoF mutation frequencies**

42,3%

1,4%
1,7%

33,8%

20,8%

**Observed LoF mutation frequencies**

45.7%

64(19,8%)                                  84(25,9%)

5(1,5%)
13(4%)

118(36,4%)                                 40(12,4%)

■ Frameshifts    ■ Deletions    ■ Insertions    ■ Stops

■ Splice mutations    ■ Start losses    ■ Stop losses

# C

# Distribution of the number of LoF mutations per gene

A *Lucifuga dentata*

B *Lucifuga gibarensis*

# A

# Eye genes

# B

# Circadian clock genes



# C

# Pigmentation genes

Probability of 19 pseudogenes

## A  Pseudogenes in cavefishes and fossorial mammals eye genes

**Cavefishes**

**Fossorial mammals**

gene
n n n
n n

Legend:
- S.anshuiensis
- C.asiatica
- C.cristata
- H.glaber

gene
n n n
n n

para 1 | 2
opn7 1|1 1|1|1
tmt 1|1 1|1|1|1
cryaa 1
parie 1 | 1
crybb 1 1|1|2
crybgx 2
crygn2 1 1
opn3 1
opn6 3|2
rrh 1
va 1
gcap5 1
opn4x 1|1
gcap7 1
cryba 1
gcap3 1
crygm5 2|1

lws 7
gnat2 6 2
grk1 1 2
grk7 5 1 7 1|4
gc2 1 1
gc3 1 2
sws1 1 5
gnb3 1
gngt2 1
rgr 9 2 1
opn4 2 1|2 1|3
gucy2f 28 11
pde6ga 1
rh1 1
pde6c 6 1 1|1
pde6b 1
pde6a 2
arr3 16 8 1
pde6h 1 1
gcap2 5 4 1 3
rpe65 1|1
rcv2 1

abca4 28
c2orf71 1
cngb3 16 1
best1 1 3
crb1 18 54 5
impg2 3
rdh5 10
rom1 12 12
rp1l1 5 4
slc24a1 1 6
rbp3 2

## B  Distribution of LoF mutations in fossorial mammals and cavefish pseudogenes

**% of genes** (y-axis, 0–100)

**Number of LoF mutations** (x-axis: 1, 2, 3, 4, 5, >5)

Data labels:
- 1: 7, 14, 2, 2
- 2: 2, 4, 1
- 3: 1, 1
- 4: 1, 1
- 5: 1, 2
- >5: 5, 2, 1