

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15

Additional Notes for

The *Sinocyclocheilus* cavefish genome provides insights into cave adaptation

Junxing Yang*†, Xiaoli Chen †, Jie Bai†, Dongming Fang†, Ying Qiu†, Wansheng Jiang†, Hui Yuan, Chao Bian, Jiang Lu, Shiyang He, Xiaofu Pan, Yaolei Zhang, Xiaoi Wang, Xinxin You, Yongsu Wang, Ying Sun, Danqing Mao, Yong Liu, Guangyi Fan, He Zhang, Xiaoyong Chen, Xinhui Zhang, Lanping Zheng, Jintu Wang, Le Cheng, Jieming Chen, Zhiqiang Ruan, Jia Li, Hui Yu, Chao Peng, Xingyu Ma, Junmin Xu, You He, Zhengfeng Xu, Pao Xu, Jian Wang, Huanming Yang, Jun Wang, Tony Whitten*, Xun Xu*, Qiong Shi*

† These authors contributed equally to this work.
* To whom correspondence should be addressed. E-mail: yangjx@mail.kiz.ac.cn (JY), Tony.Whitten@fauna-flora.org (TW), xuxun@genomics.cn (XX), or shiqiong@genomics.cn (QS)

Contents

16		
17		
18	Note S1	3
19	1. Organism background, genome sequencing and assembly	3
20	1.1 Organism background	3
21	1.2 Sample background and sequencing	3
22	1.3 Genome size estimation by k-mer analysis and genome assembly	3
23	1.4 Assessment of genome assemblies	4
24	Note S2	6
25	2. Genome annotation	6
26	2.1 Repetitive sequence detection	6
27	2.2 Protein-coding gene annotation	6
28	2.3 Non-coding RNA annotation	7
29	Note S3	9
30	3. Transcriptome analysis	9
31	3.1 RNA isolation, and RNA-Seq method	9
32	3.2 RNA-Seq data processing	9
33	3.3 Alignment and <i>de novo</i> assembly of the transcriptomes	9
34	Note S4	10
35	4. Evolutionary analysis of <i>Sinocyclocheilus</i> genomes	10
36	4.1 Gene family cluster	10
37	4.2 Phylogenetic analysis	10
38	4.3 Divergence time estimation	11
39	4.4 Demographic history	11
40	4.5 Gene family contraction and expansion	11
41	4.6 Genes with accelerated evolutionary rates	12
42	4.7 Evolution of Hox clusters	13
43	4.8 Loss of Sa-specific gene families	14
44	Note S5	16
45	5. Morphological comparison analysis	16
46	5.1 Paraffin section of eyes of the three <i>Sinocyclocheilus</i> species	16
47	5.2 Immunocytochemistry of taste buds	16
48	5.3 Anatomy of gas bladder, absolute fecundity and other measurements	17
49	5.4 Synchrotron X-ray microtomography of the saccular otolith	17
50	Note S6	18
51	6. Cave adaption analysis	18
52	6.1 Vision analysis	18
53	6.2 Pigmentation	19
54	6.3 Scale development	20
55	6.4 Hearing	20
56	6.5 Immune response	21
57	6.6 Circadian rhythm	21
58	6.7 Sense of taste	22
59	References (1-85)	23
60		

61 **Note S1**

62 **1. Organism background, genome sequencing and assembly**

63 1.1 Organism background

64 The cyprinid genus *Sinocyclocheilus* (Cypriniformes: Cyprinidae) is endemic to the massive
65 southwestern karst area which abuts the eastern Qinghai-Tibetan Plateau of China, covering the
66 Yunnan-Guizhou Plateau and the surrounding region (including east Yunnan Province,
67 south-central Guizhou Province, and northwest Guangxi Province). There are about 50 valid
68 species distributed within a relatively narrow area (about 270,000 km²), making it the most
69 species-rich cyprinid genus in China [1]. *Sinocyclocheilus* is well known for its marked
70 adaptations to the subterranean environment. Its high species diversity and phenotypic variation
71 make it an emerging model for studying the evolution of cavefishes.

72 Three adult female fishes (Sg: *S. grahami*, Sr: *S. rhinoceros* and Sa: *S. anshuiensis*) were
73 selected for whole genome shotgun sequencing in this study, which were collected in Lake
74 Dianchi, Yunnan; Luoping, Yunnan; and Lingyun, Guangxi respectively; representing
75 surface-dwelling, semi-cave-dwelling and cave-restricted species (Figure S1). The research
76 protocol and treatment of experimental animals were reviewed and approved by the internal
77 review board of the Kunming Institute of Zoology, Chinese Academy of Sciences (approval ID:
78 SYDW-2014020).

79

80 1.2 Sample background and sequencing

81 High-quality genomic DNA was extracted from the muscle tissues using Puregene Tissue Core
82 Kit A (Qiagen, Maryland, USA) for construction of libraries with different inserted sizes (250
83 bp~20 kb) (Table S1). In total, twenty-five paired-end libraries (11 for Sg, 7 for Sr and 7 for Sa,
84 respectively) were generated with the Illumina standard operating procedure. Paired-end
85 sequencing was performed on Illumina Hiseq2000 for each library.

86

87 1.3 Genome size estimation by k-mer analysis and genome assembly

88 There were 293.70, 174.02 and 188.28 Gb raw data obtained for the Sg, Sr and Sa fishes,
89 respectively (Table S1). Artificial and low-quality reads were filtered first and then 88.05, 48.31
90 and 60.54-fold coverage of *Sinocyclocheilus* genomes were used for assembly (Table S2). We
91 corrected the sequencing errors with 17-mer frequency lower than four with a method described
92 in a previous study [2]. In addition, we estimated genome sizes of the three *Sinocyclocheilus*
93 species using the 17-mer depth frequency distribution method: G (Genome size) =

94 K-mer_num/Peak_depth. The estimated genome sizes are 1.79, 1.89 and 1.81 Gb, respectively
95 (Figure S3, Table S4). We subsequently assembled the filtered reads into contigs and scaffolds to
96 build the genomes using SOAPdenovo [3] and fulfilled the gaps with GapCloser [3]. Our
97 assembly strategy allowed us to obtain a 1.75-Gb genome for Sg with a scaffold N50=1.16 Mb,
98 a 1.73-Gb genome for Sr with a scaffold N50=896.4 kb and a 1.68- Gb genome for Sa with an
99 scaffold N50=1.25 Mb (Table S3).

100 To assess the quality of our assemblies, we performed three different evaluations. First, we
101 calculated the GC content of these three fishes (37.3% for Sg, 37.2% for Sr, and 37.3% for Sa),
102 which were similar to those in zebrafish (*Danio rerio*) and medaka (*Oryzias latipes*) (Figure S3).
103 Second, the GC depth distributions were relatively concentrated with a GC content above 20x
104 (Figure S4). Third, we aligned the sequencing reads to the assembled genomes and determined
105 the sequencing depth distribution, which showed that the proportion of sequencing was close to
106 their sequencing depth, and depths lower than 10x were under 5% (Figure S5) [4].

107

108 1.4 Assessment of genome assemblies

109 1.4.1 Evaluation of the Sg assembly with *de novo* assembled fosmids

110 Three fosmid sequences were also used to check the genome assembly of Sg. Single-end
111 Sanger sequencing was performed, and the average length of generated reads was about 650 bp.
112 After filtering, about 7.5 folds of coverage for each fosmid were used for assembly.

113 The fosmids were aligned against the scaffolds using BLAT with default parameters. The
114 alignment blocks were then linked on fosmid sequences by in-house Perl scripts. Paired-end reads
115 of short inserted size (< 1000 bp) of Sg were mapped on the fosmid sequences using SOAP
116 (version 2.21). SOAPcoverage (SOAP software package) was applied to calculate sequencing
117 depth for each fosmid sequence with a non-overlapping 100-bp window.

118 As a result, each fosmid (average length=35 kb) was aligned to only one scaffold (Figure S6),
119 and up to 93.6% of the total fosmid regions (Table S5) were well covered by the assembled
120 contigs.

121

122 1.4.2 Assessment of genome assembly coverage

123 The assessment of genome assembly coverage was carried out using CEGMA program
124 (version 2.3) (<http://korflab.ucdavis.edu/datasets/cegma/>) [5, 6], which combin í tblastn
125 (blast-2.2.25), genewise (wise2.2.3), hmmer (hmmer-3.0), and geneid (version 1.4). By
126 comparing the completeness of predicted core eukaryotic genes (CEGs) with 248 highly
127 conserved CEGs in the database, we could evaluate the coverage of our assembled genomes.

128 The results indicated that the ratio of complete CEGs coverage were 79.84%, 81.05% and
129 85.48% respectively, and the ratios of partial coverage of the three species were more than 96%
130 (Table S9).

131

132 1.4.3 Pairwise whole genome alignment

133 Pairwise whole genome alignment among Sg, Sr and Sa was carried out using LASTZ, with the
134 parameters setting as: C=2, T=2, H=2000, Y=3400, L=6000, and K=2200. Subsequently the
135 Chain/Net package was used for post treatment. The genomes were masked with RepeatMasker
136 repeats at “-s” setting and TRF tandem repeats of period ≤ 12 (Table S6). The aligned length
137 between each two was about 760 Mb (Table S7).

138

139 1.4.4 Multiple whole genome alignment

140 The 3-way whole genome multiple alignment of Sg, Sr and Sa was generated using multiz
141 following the topology of species tree. The Sa genome was set to be the reference, and for input
142 pairwise alignments we carried out in-house generation of the Sg versus Sr and Sg versus Sa
143 alignments (Figure S7).

144

145 1.4.5 Detection of segmental duplication in *Sinocyclocheilus* genomes

146 We analyzed the Sg, Sr and Sa genome assemblies using Whole Genome Analysis Comparison
147 (WGAC) approaches that were designed to detect genomic duplicates >1-kb length and >90%
148 sequence identity.

149 We applied LASTZ designed for long genomic sequence alignments within Sg, Sr and Sa
150 assemblies. First, self-versus-self LASTZ alignment was performed using the repeat-masked
151 genome sequences, and then chaining of well-ordered neighboring alignments. We obtained the
152 seeding segmental duplications (non-repeat alignment length >500 bp and overall identity >85%),
153 and then reintroduced the masked repeat regions to perform optimal global alignment to refine the
154 alignment identity and define the boundaries of segmental duplications more accurately. The
155 resulting alignments that extended to >1-kb length and had >90% sequence identity were deemed
156 segmental duplications (Table S8).

157 From the results, we observed that segmental duplications in the *Sinocyclocheilus* genomes
158 were more than those in other fishes. However, we are not sure whether the results were due to
159 chromosomal rearrangements or chromosome doubling.

160

161 **Note S2**

162 **2. Genome annotation**

163 2.1 Repetitive sequence detection

164 Almost every assembled genome is composed of repeated sequences. There are two main types
165 of repeats in the genomes, including tandem repeats and interspersed repeats (such as LINEs and
166 SINEs). We identified repetitive sequences or transposable elements (TEs) of three fishes based
167 on combination of homology and *de novo* prediction methods using a variety of softwares [7].

168 We performed homology-based prediction by running the RepeatMasker (version 3.3.0) [8]
169 with the default parameters against the Repbase [9], and using TE library (version 16.10) as well
170 as RepeatProteinMask also with the default parameters at the DNA and protein levels to identify
171 repeat sequences. For *de-novo* prediction of TEs, we searched the denovo-prediction-build repeat
172 library using RepeatModeler (version 1.0.5) and generated TE results with classification
173 information for each repeat family by running RepeatMasker on the genome sequences
174 subsequently. Tandem repeats were also searched using the Tandem Repeats Finder (version 4.04)
175 [10] with parameters “Match=2, Mismatch=7, Delta=7, PM=80, PI=10, Minscore=50, and
176 MaxPeriod=2000”. We generated the final results (Tables S10 to S12) with integrating all of the
177 repeat annotation results using an in-house perl program. The sequence divergence rate was also
178 calculated for each family of TEs (Figure S8).

179

180 2.2 Protein-coding gene annotation

181 The protein coding genes were obtained using combination of *de novo* method,
182 homology-based gene prediction and RNA-seq data. All predicted gene evidence was integrated
183 by GLEAN [11] to get non-redundant data [12].

184 For homology-based gene prediction, protein coding sequences of zebrafish, three-spined
185 stickleback (*Gasterosteus aculeatus*), medaka, green spotted puffer (*Tetraodon nigroviridis*) and
186 human (*Homo sapiens*) from Ensembl database (release 64) were used for analysis. The longest
187 transcripts were chosen to represent the alternative splicing variants of genes [4]. Firstly, we
188 aligned our assembled genomes to those protein sequences using TBLASTN (blast-2.2.23) [13]
189 with E-value < 1e-5 and dealt with the alignment results by SOALR. Then GeneWise (version
190 2.2.0) [14] was used to predict the precise transcript structure and filter low-quality results with
191 cds <150 bp after clustering the gene set. Therefore, the homology-based gene sets were obtained.

192 Before performing *de novo* prediction, we masked all the repetitive sequences in the assembled
193 genomes. Then we carried out annotation using Augustus [15], SNAP [16] and GlimmerHMM

194 [17] with gene model parameters trained by 800 completed genes selected from homolog
195 alignment, and finally obtained the intersection of results from the three methods.

196 In addition, RNA-seq data from different tissues including eye, skin, liver and gonad of three
197 fishes were also used to verify the gene sets. First, we mapped transcriptome data to the genome
198 using Tophat (v.2.0.4) [18], and then got transcriptome-based gene structures by cufflinks (v.2.0.0)
199 [19]. Lastly, all gene structures obtained by the three approaches were combined using GLEAN to
200 filter genes with low confidence level.

201 Using the above approaches, we obtained the final gene sets of the three *Sinocyclocheilus*
202 species (Sg: 42,109 genes, Sr: 40,333 genes and Sa: 40,470 genes) (Table S13). CEGMA was
203 used again to assess the coverage rate between KOG (EuKaryotic Orthologous Groups) genes
204 predicted by CEGMA and the three gene sets in this paper (Table S14). The results indicate that
205 the predicted gene sets cover more than 98% of KOGs at least. Meanwhile, gene numbers in each
206 set are nearly as twice as the number of zebrafish (26,206), which supports that the
207 *Sinocyclocheilus* lineages are tetraploid. The distributions in mRNA, CDS, exon and intron
208 lengths of protein coding genes in each genome with a comparison to zebrafish are also
209 summarized (Figure S9). Additionally, we identified co-linearity blocks between zebrafish and Sg
210 genomes regarding adjacent synteny genes as a block. With more than half of blocks have 2:1
211 relationship between Sg and zebrafish (Figure S10), we infer that there was another genome
212 duplication in *Sinocyclocheilus* (maybe at an unknown ancient node) which is called the 4R
213 duplication event.

214 In order to further understand the species' biological adaptation, we explored their gene
215 functions. First, we used our coding proteins blast against databases including PRINTS [20],
216 PROSITE [21], Pfam [22], ProDom [23], SMART [24] and PANTHER [25] by InterProScan [26]
217 to annotate motifs and domains. Meanwhile, GO (Gene Ontology) [27] items for each gene were
218 gained from the InterProScan result. Then we also searched KEGG (Kyoto Encyclopedia of
219 Genes and Genomes) databases [28] (release 58) to find the pathways in which the gene might be
220 involved, as well as Swissprot (version 2012.03) and TrEMBL (version 2012.03) to confirm the
221 gene symbols. The final results are listed in Table S15.

222

223 2.3 Non-coding RNA annotation

224 A non-coding RNA (ncRNA) is any RNA molecule that is not translated into a protein. Four
225 types of ncRNA, including micro RNA (miRNA), transfer RNA (tRNA), ribosomal RNA (rRNA),
226 and small nuclear RNA (snRNA) were annotated in our analysis. There are four types of rRNA
227 composed of 5S, 5.8S, 18S and 28S in eukaryotes, which are quite conservative. Hence, we

228 searched rRNAs by blasting human rRNA database. As for siRNAs and snRNAs, we blasted
229 against Rfam [29, 30] database which divided ncRNAs into families based on evolution from a
230 common ancestor. Finally, we used tRNAscan-SE [31] to search reliable tRNAs in the genomes.
231 The detailed information is shown in [Tables S16 to S17](#).

232

233

234 **Note S3**

235 **3. Transcriptome analysis**

236 3.1 RNA isolation, and RNA-Seq method

237 RNA was isolated for sequencing from four tissues (eye, skin, liver, and ovary) of Sg, Sr and
238 Sa respectively. RNA-seq libraries were constructed using the Illumina mRNA-Seq Prep Kit, then
239 the libraries were sequenced using Illumina sequencing platform and the 90-bp paired-end
240 sequencing module. About 4-Gb clean data were generated in each library.

241

242 3.2 RNA-Seq data processing

243 There are some adaptor sequences and/or low-quality reads presented in the raw reads. In order
244 to obtain high-quality reads, we applied an in-house C++ program to filter our raw reads. The
245 procedure included the following steps:

246 Remove reads with adaptor sequences.

247 Remove reads that contained > 10% ambiguous base calls (Ns).

248 Remove low-quality reads that contained > 40% low-quality base calls (quality score ≤ 5).

249 Remove read pairs with read 1 and read 2 overlapping by ≥ 10 bp and 100% identical read
250 pairs.

251

252 3.3 Alignment and *de novo* assembly of the transcriptomes

253 All of the clean RNA-Seq reads were mapped onto the corresponding reference genomes (Sg,
254 Sr and Sa) using TopHat. According to the mapped results, transcripts were constructed using
255 cufflinks. The *de novo* transcriptomes of the four tissues were assembled by Trinity with filtered
256 reads from each tissue separately into contigs and scaffolds. Trinity contains Inchworm, Chrysalis
257 and Butterfly, which were employed sequentially to process large volumes of RNA-seq reads.

258

259

260 **Note S4**

261 **4. Evolutionary analysis of *Sinocyclocheilus* genomes**

262 4.1 Gene family cluster

263 Orthologs refer to homologous sequences derived by a speciation event from a single ancestral
264 sequence in the last common ancestor of the examined species. We defined gene families using
265 TreeFam (<http://www.treefam.org>) among the three *Sinocyclocheilus* species (Sg, Sr and Sa) and
266 seven other vertebrates, including fugu, green-spotted puffer, three-spined stickleback, Atlantic
267 cod, medaka, zebrafish, and human. First, we ran all-vs-all blast with the e-value of $1e-7$ using ten
268 species's protein sequences to rapidly find an initial list of possible matches. Then, we conjoined
269 the blast alignments by Solar, removed genes with aligned ratio < 0.33 , converted the bit score to
270 percent score, and grouped the genes by hcluster_sg (<https://pypi.python.org/pypi/hcluster>,
271 version 0.5.0). Eventually, 17,883 gene families and 210 single-copy gene families were
272 generated. The numbers of orthologous genes across the ten species were counted ([Figure S13](#),
273 [Table S19](#)) and plotted in a Venn diagram ([Figure S14](#)).

274

275 4.2 Phylogenetic analysis

276 We constructed phylogenetic trees using the gene sets of Sg, Sr and Sa, together with seven
277 other published vertebrate genomes (fugu, green-spotted pufferfish, three-spined stickleback,
278 Atlantic cod, medaka, zebrafish and human). Since *Sinocyclocheilus* are tetraploid, the
279 single-copy families were rather fewer than in other vertebrates. We extracted these genes with
280 two-copy families in *Sinocyclocheilus* but one-copy in the other seven vertebrates, and then
281 removed the shorter one of the two copies in *Sinocyclocheilus*, hence the remaining longer genes
282 were presented as a group of single-copy families. Together with the original 210 single-copy
283 families, we obtained 3,181 single-copy families. Subsequently, we extracted amino acid
284 sequences from the 3,181 single-copy genes to concatenate one super-gene, and constructed
285 phylogenetic trees using maximum likelihood (ML) method in PhyML [32, 33] and Bayes
286 inference method in MrBayes [34]. Different datasets from whole coding sequences, 1st, 2nd, 1st &
287 2nd, 3rd and 4d codon positions were used to reconstruct phylogenetic tree respectively. All results
288 revealed that Sg and Sa clustered together and then grouped with Sr with high bootstrap values.

289 Meanwhile, six mitochondrial gene (*COX1*, *COX2*, *COX3*, *ND3*, *ND4* and *ND6*) sequences
290 were used to reconstruct the phylogenetic trees with the same pipeline and parameters as
291 described above [35], which also supported the same phylogenetic relationships. These topologies
292 of the species tree are shown in [Figure S15](#) and [S16](#).

293

294 4.3 Divergence time estimation

295 Based on the phylogenetic tree and sequences derived from CDS, 1st and 2nd and 4d codon
296 positions, and using five fossil records (Takifugu-Tetraodon, 32.25~56 Ma; Stickleback-Takifugu,
297 Tetraodon, 96.9~150.9 Ma; Medaka-Stickleback, Takifugu, Tetraodon, 96.9~150.9 Ma;
298 Zebrafish-Medaka, Stickleback, Takifugu, Tetraodon, 149.85~165.2 Ma; Fishes-toad, bird,
299 mammal, 416~421.75 Ma) [36] as calibrating time points, we estimated the species divergence
300 time using the mcmctree [37] in PAML [38] with default values in parameter settings.
301 Furthermore, we also recalculated the divergence time using the multidivtime [39, 40] program.

302 The generated chronogram using mcmctree in PAML based on 4d codon positions revealed that
303 Sr emerged first from other two species at 26.3 Ma, and Sg divided from Sa at 17.5 Ma. We also
304 obtained similar results when using other methods (Multidivtime, CDS sequences) and datasets
305 (1st & 2nd) to estimate divergence times, which are shown in [Figure S17](#).

306

307 4.4 Demographic history

308 The distribution of time to TMRCA (the most recent common ancestor) between two alleles in
309 an individual can be related to the history of population size fluctuation. The population size
310 histories of Sg, Sr and Sa were inferred using the pairwise sequentially Markovian coalescent
311 (PSMC) model [41] on heterozygous sites with the generation time ($g=1$ year) (according to
312 artificial breeding of Sg) and the mutation rate ($m=3.51 \times 10^{-9}$ per year per nucleotide) [42].
313 Reconstructed population history was plotted for Sg, Sr and Sa separately using gnuplot4.4 [43].

314

315 4.5 Gene family contraction and expansion

316 According to the above-mentioned phylogenetic trees and divergence time, we explored gene
317 families that underwent a significant change in *Sinocyclocheilus* lineages. After filtering the
318 families that have only one copy member or those have more than 200 copies in one species while
319 other species had less than 2 copies, we performed *Sinocyclocheilus* lineage-specific expansion
320 and contraction analysis using the CAFÉ program [44]. Based on a random birth and death model
321 [45], a global parameter λ was estimated using maximum likelihood. Then, gene family numbers
322 in each species and their ancestor were estimated. Comparing each branch and their ancestor
323 branch, a conditional p-value was calculated and families with p-value less than 0.05 were
324 marked as a significant change [46], which means it underwent contraction or expansion in the
325 process of evolution ([Figure 2a](#)).

326 After that, the families that significantly change in size were subjected to GO/KEGG/IPR

327 enrichment analyses along each *Sinocyclocheilus* lineage, respectively (Table S24).

328 Among those gene families undergoing expansion and contraction, a contractive family *OPN5*
329 in Sa (GO: 0007601: visual perception) caught our attention, because it may be associated with
330 eye degeneration for cave adaption. Opsins are members of the guanine nucleotide-binding
331 protein-coupled receptor superfamily. *OPN5*, which is expressed in the eye, brain, testes, and
332 spinal cord, belongs to the seven-exon subfamily of mammalian opsin genes that includes
333 peropsin (RRH) and retinal G protein coupled receptor (RGR) and encodes a protein with
334 photoisomerase activity [47]. The olfactory receptor family is the largest one in the genome,
335 responsible for the recognition and G protein-mediated transduction of odorant signals. We found
336 the *OR* family (GO:0004984 olfactory receptor activity) has undergone a large expansion,
337 especially *OR10*, *OR13*, *OR2*, *OR51*, *OR52*, *OR1* and *OR6*. There were some other families
338 associated with immune response, such as *DPBI* family (GO:0042613 MHC class II protein
339 complex) which have experienced contraction, while *HLA-B* family (GO:0042612 MHC class I
340 protein complex) has experienced expansion.

341 More details are shown in Table S29 in additional file 2 (Table S29-genefam.expan.contract.list.
342 xlsx)

343

344 4.6 Genes with accelerated evolutionary rates

345 Positive Darwinian selection at the DNA sequence level has been tested by estimating the ratio
346 (ω) of nonsynonymous nucleotide substitutions (dN) to synonymous nucleotide substitutions (dS)
347 between ortholog genes [48]. We performed positive selection analysis based on the single-copy
348 genes that we used in the phylogenetic analysis. After filtering families with an alignment ratio of
349 less than 0.4, we set the three *Sinocyclocheilus* species as the foreground branch respectively,
350 while other species in our phylogenetic tree were set as the background branch. A branch-site
351 model [49] was used to calculate different ω in phylogenetic tree, such as ω_0 , ω_2 and ω_1
352 representing the whole tree, foreground branch and background branch respectively [35]. Then a
353 p-value test was used to check whether ω_2 was significantly greater than ω_1 and ω_0 . The genes
354 with a p-value less than 0.05 implied that they may have undergone positive selection in the
355 *Sinocyclocheilus* lineages.

356 After obtaining the *Sinocyclocheilus* positive selection gene (PSG) list, we converted it to the
357 corresponding human orthologs as the input against a background of human genes [48] using
358 DAVID Functional Annotation [50] tool. The annotation and enrichment result files are listed in
359 Table S30 in additional file 3 (Table S30-positive.selection.list.xlsx).

360 In the Sa branch, the following genes should be given special attention: *DHX30*, *Rab3gap1*,

361 *BBS4*, *Gpr143* and *galE*. *DHX30* encodes a mitochondrial nucleotide protein which is identified
362 as a component of a transcriptional repressor complex that functions in retinal development.
363 Defects in *Rab3gap1* are the cause of Warburg micro syndrome 1 [51, 52], which is a severe
364 autosomal recessive disorder characterized by developmental abnormalities of the eye and central
365 nervous system. Defects in *BBS4* are the cause of Bardet-Biedl syndrome type 4 (BBS4) that is
366 usually characterized by severe pigmentary retinopathy, early onset obesity, polydactyly,
367 hypogenitalism, renal malformation and mental retardation [53]. *Gpr143* encodes a protein that
368 binds to heterotrimeric G proteins and is targeted to melanosomes in pigment cells. Mutations in
369 this gene cause ocular albinism type 1 [54], especially associated with visual impairment.
370 Mutations in *galE* result in epimerase-deficiency galactosemia, also referred to as galactosemia
371 type 3, a disease characterized by liver damage, early-onset cataracts, deafness and mental
372 retardation [55]. All the above genes may be connected to physiological and morphological
373 changes in cavefishes.

374 In the Sr branch, we also found *DHX30* and *Rab3gap1* underwent positive selection.
375 Accordingly, we found *DGKI*, a member of type IV diacylglycerol kinase subfamily, may play a
376 crucial role in the production of phosphatidic acid in the retina or in recessive forms of retinal
377 degeneration (from NCBI, no reference). Mutations in *HPS1* may play a role in organelle
378 biogenesis associated with melanosomes, platelet dense granules and lysosomes, and are
379 associated with Hermansky-Pudlak syndrome type 1 which characterized by oculocutaneous
380 albinism, bleeding and lysosomal storage defects [56]. *GIPC3* has also undergone positive
381 selection; its mutations are associated with autosomal recessive deafness [57].

382

383 4.7 Evolution of Hox clusters

384 With the highly conserved sequences and synteny genomic architecture, the Hox clusters are
385 usually used as evidence of whole genome duplication [58]. More importantly, the order of Hox
386 genes along the chromosome reflects the order in which they act along the body. Teleost fishes
387 are the pinnacle of Hox cluster evolution, with at least seven Hox clusters and 49 genes in
388 zebrafish [58].

389 To define Hox genes in the three *Sinocyclocheilus* genomes, the Hox genes of zebrafish were
390 downloaded from ensembl, then reciprocal best ortholog analysis was used with an e-value cutoff
391 of 1e-5 and filtered with less than 80% of alignment and 75% overall amino acid identity.

392 The gene numbers for the *Sinocyclocheilus* species were nearly twice of those in zebrafish
393 (Table S18). Remarkably, the order of Hox genes along the scaffolds in *Sinocyclocheilus* reflects
394 the same order as in zebrafish, and the diagrams compared syntenic relationships with the Hox of

395 zebrafish linkage groups. The multiple copies of Hox genes indicated that *Sinocyclocheilus* fishes
396 are indeed tetraploids.

397 Hox genes define patterns of development in vertebrate limbs. For each Hox cluster composed
398 of Anterior (group 1,2), Group3 (group 3), Central (group 4~8) and Posterior (9~13) classes of
399 Hox genes existed before the divergence of the protostome and deuterstome lineages. From the
400 results (Figure S11), we found that the conserved non-coding sequences encompass a surprisingly
401 large part of the clusters, especially in the *HoxBa* and *HoxCa* clusters.

402 Group 1 appeared in *HoxA*, *HoxB* and *HoxC* clusters of zebrafish and *Sinocyclocheilus*
403 genomes, therefore it was used to construct a similarity network. The results (Figures S12A to
404 S12B) showed that each of these three clusters of *Sinocyclocheilus* was orthologous to one of the
405 zebrafish clusters. Hence, the duplication events of *Sinocyclocheilus* that produced the clusters
406 occurred at least before the divergence of *Sinocyclocheilus* lineages at about 26.3 million years
407 ago (Figure 2a).

408

409 4.8 Loss of Sa-specific gene families

410 In order to identify Sa-specific gene family loss, we extracted gene families that have no
411 member in Sa while being larger than zero in the other nine species. The lost gene family list is
412 included in Table S23. Many genes, such as *LMCD1*, *TEM-7*, *PDCD5*, *BCL6*, *F7*, *CA14*,
413 *CYP2U1* and *Creb3l4*, have lost at least one copy in Sa. Previous studies have shown that
414 *Creb3l4* can regulate the expression of genes required for germ cell survival but is insufficient to
415 disrupt the normal fertility in mice [59]. We infer that the loss of *Creb3l4* may have some
416 relationship with its low fecundity. The inner mitochondrial membrane protein Mpv17 in
417 podocytes is essential for the maintenance of mitochondrial homeostasis and protects podocytes
418 against oxidative stress-induced injury both in vitro and in vivo [60]. *LMCD1*, *TEM-7* (tumor
419 endothelial marker 7), *PDCD5* were associated with tumor metastasis [61], antiaging [62] and
420 enhanced apoptosis [63] respectively. *BCL6* may function in a 'hit-and-run' role in
421 lymphomagenesis [64]. *F7* (EC 3.4.21.21, blood-coagulation factor VIIa, activated blood
422 coagulation factor VII, formerly known as proconvertin) is one of the proteins that causes blood
423 to clot in the coagulation cascade. *CA14* is a Carbonic anhydrase (CA) that catalyzes the
424 reversible hydration of carbon dioxide. It participates in a variety of biological processes,
425 including respiration, calcification, acid-base balance, bone resorption, and the formation of
426 aqueous humor, cerebrospinal fluid, saliva, and gastric acid. It has shown extensive diversity in
427 tissue distribution and in their subcellular localization [65]. *CYP2U1* (cytochrome P450, family 2,
428 subfamily U, polypeptide 1) metabolized arachidonic acid, docosahexaenoic acid (DHA), and

429 other long chain fatty acids, suggesting a potential role in brain and immune functions [66].
430 *TSC2* in Sa is under positive selection. Its gene product is believed to be a tumor suppressor
431 and is able to stimulate specific GTPases, which is consistent with the gene loss of *LMCD1*,
432 *TEM-7*, *PDCD5* for tumor suppression.
433

434 **Note S5**

435 **5. Morphological comparison analysis**

436 5.1 Paraffin section of eyes of the three *Sinocyclocheilus* species

437 Histological studies on eyes were performed according to the following procedures: (1)
438 Fixation: Separated eye regions were fixed in 4% paraformaldehyde at 4°C overnight. (2)
439 Dehydration: Eye tissue was dehydrated through a series of alcohols, including 30%, 50%, 70%,
440 85%, two 95%, and three 100% steps, 20 min per step. They were transferred into
441 ethanol/TO-reagent (1:1) twice, and then moved into TO reagent twice more. (3) Embedment:
442 Paraffin/TO-reagent (1:1) was added for one hour (62°C), followed by 100% paraffin twice in two
443 hours (62°C). The tissue was embedded into melting and solidifying paraffin. (4) Sectioning:
444 Tissues were sectioned at 5–15 µm thickness in the microtome blocks, and then spread the
445 sections into glass slides. (5) Dewaxing the Slides and Rehydration of Tissues: Dewaxed slides
446 were placed into two TO containers for ten minutes each, and then the rehydrated slides were
447 placed into graded ethanol solutions: 100%, 95%, 90%, 80%, and 70% for 3-5 minutes each. (6)
448 Staining: Sections were stained using Hematoxyli-Eosin (HE). (7) Photography: The sections
449 were checked on the slides and photographed under a microscope.

450

451 5.2 Immunocytochemistry of taste buds

452 Immunocytochemistry analysis of taste buds (checking only the separated upper and lower
453 jaws, hereafter called tissues) was carried out according to the following steps: (1) Fixation:
454 Tissues were fixed in 4% paraformaldehyde in PBS overnight at 4°C. (2) Immunocytochemistry:
455 The tissues were washed thoroughly in PBS and then permeabilized overnight at 4°C in PBS
456 containing 4% Triton X-100. (3) The tissues were washed thrice more the following day for 15
457 min each with PBS at room temperature. They were incubated in primary antibody (Rabbit
458 against calretinin, labeling entire receptor cells within taste buds) and occasionally agitated gently
459 in a dark chamber for 5-6 days at 4°C. This antibody was added together in PBS with 0.5% Triton
460 X-100 to yield final dilutions of 1:1500. (4) The tissues were rinsed three times for 30 min in each
461 with PBS. (5) The tissues were incubated in secondary antibodies (a final dilution of 1:500).
462 Incubation in secondary antibody similarly occurred with gentle agitation in the dark chamber for
463 further 2-3 days at 4°C. (6) The tissues were rinsed three times for 45 min each with PBS. (7)
464 Negative controls were conducted in tissues that were incubated in secondary antibody without

465 first being labelled with primary antibody. (8) The tissues were viewed using a microscope
466 equipped for epifluorescence.

467

468 5.3 Anatomy of gas bladder, absolute fecundity and other measurements

469 Samples for anatomy and measurements were from specimens immersed in 90% ethanol. We
470 used dissecting scissors and tweezers to separate the gas bladder from other viscera and then took
471 photos and measured them. Absolute fecundity was counted from the number of eggs in
472 individuals at maturity stages IV (mature) or V (ripe). The counting of lateral line scales was
473 conducted under a stereo microscope.

474

475 5.4 Synchrotron X-ray microtomography of the saccular otolith

476 The synchrotron X-ray microtomography experiments were performed at BL13W1 beamline of
477 Shanghai Synchrotron Radiation Facility. A specimen was held in a vertical tube mounted on a
478 sample stage. If the specimen was larger than the X-ray beam, the sample stage was moved
479 vertically for continuous scans. Specimens were imaged with 13 keV monochromatic X-ray at 9
480 μm resolution for larger specimens, and at 3.7 μm for juvenile specimens. The slices were
481 reconstructed using the FBP algorithm, and 3D renderings were created and manipulated in
482 VGStudio 2.1 software.

483

484 **Note S6**

485 **6. Cave adaption analysis**

486 6.1 Vision analysis

487 Protein sequences of 140 zebrafish vision-related genes were downloaded from the present
488 Ensembl database and aligned to our gene sets using blastp with E-value <1e-5. The homolog
489 genes of the three *Sinocyclocheilus* species were determined according to the identity with
490 zebrafish's category and the copy numbers were counted (Tables S25A to S25E). Multiple protein
491 sequence alignments of homolog genes were performed with MUSCLE (version 3.8.31) [67]. The
492 results demonstrated that Sg *rcv-1* gene was premature and Sa *rom1b* gene has a two-amino-acid
493 insertion and the inserted site is located in the Tetraspanin family conserved domain (also
494 testified by PCR) [NCBI accession: pfam00335] (Figure S23). The sequencing depth near the
495 insert site was found with IGV (Integrative Genomics Viewer) software (version 2.3.32) (Figure
496 S23).

497 A similarity cluster of crystalline genes of zebrafish and three *Sinocyclocheilus* species is
498 shown in Figure S28. The *Sinocyclocheilus* genes were classified according to the zebrafish gene
499 nomenclature. Among 60 crystallin genes, 29 zebrafish genes were defined as outgroup compared
500 with *Sinocyclocheilus* genes, however, 4 zebrafish genes cannot be distinguished from those of
501 the *Sinocyclocheilus* species. *Crygm2d* gene in the *Sinocyclocheilus* fishes cannot be classified
502 into any particular category, although it was the closest to *Crygm2d12* of zebrafish (Figure S28).

503 Development and maintenance of photoreceptors requires a series of transcriptional factors
504 including *Crx* (cone-rod homeobox), *Nrl* (neural retina leucine zipper protein), *Otx2*
505 (orthodenticle homolog-2), *Otx5* (orthodenticle homolog-5), *Nr2e3* (nuclear receptor subfamily 2
506 group E member 3), *Thrb* (thyroid hormone receptor beta), *Gucal1a* (guanylate cyclase activator
507 1A), *Gnat1*, *Gnat2*, *Irbp* (interphotoreceptor retinoid-binding protein), *Rorb* (RAR related orphan
508 receptor beta) [68]. Our eye transcriptome analysis of the three *Sinocyclocheilus* species showed
509 that the expression levels of transcriptional factors were Sg>Sr>Sa (Tables S25A to S25E, S26A
510 to S26D, Figures S25A to 25B). Ten transcriptional factors (*Crx*, *Nrl*, *Otx2*, *Otx5*, *Nr2e3*, *Gucal1A*,
511 *Gnat1*, *Gnat2*, *Rorb*, *Rx3*) were significantly down-regulated in Sa compared with Sg and Sr. On
512 the contrary, some intriguing non-vision direct related genes showed different expression patterns
513 in the three species. For example, the expression (RPKM value) of *Irf6* (interferon regulatory
514 factor 6) in Sa was 158.6, which was significantly higher than Sg (2.2) and Sr (0). *Irf6* is related
515 to the formation of connective tissue and play a critical role in keratinocyte development [69].
516 Tubulins are a small family of globular proteins. The most common members of the tubulin

517 family are α -tubulin and β -tubulin which make up microtubules. The expression of tubulin alpha
518 chain in Sa was 18.9 while there was no expression in Sg and Sr. These genes may be involved in
519 formation of connective tissue and degeneration of eyes in the cave-restricted Sa.

520 Meanwhile, the expression of *peripherin2* (photoreceptor outer segment membrane
521 glycoprotein 2) was detected in Sg (90.1) and Sr (1.5) while there was no expression in Sa.
522 Secreted extracellular proteins are often glycosylated to be functional and we speculated that the
523 function of photoreceptor outer segment membrane may decline in Sa. Oat (Ornithine
524 aminotransferase) is an enzyme involved in the ultimate formation of proline from ornithine. A
525 deficiency of this enzyme causes Oat deficiency, also known as gyrate atrophy of choroid and
526 retina. The presenting symptom of Oat deficiency is myopia which progresses to night blindness
527 [70]. No expression of *Oat* in Sa and Sr is consistent with the partial or entire dark environment.
528 *Pde6h-b* is a gene that encode retinal rod rhodopsin-sensitive cGMP 3',5'-cyclic
529 phosphodiesterase subunit gamma which is related to photosensitivity in rod photoreceptor cells
530 [71]. The expression of *Pde6h-b* in Sa (6.3) was significantly lower than those in Sg (4174) and
531 Sr (626). These observations may suggest different developmental genetic mechanisms of retinal
532 degeneration occurred in the cavefish Sa.

533

534 6.2 Pigmentation

535 To define orthology genes related to melanogenesis in the three *Sinocyclocheilus* genomes, a
536 blastp analysis was employed with the e-value cutoff of $1e-5$ using Sg, Sr, Sa coding proteins,
537 respectively. When setting human pigment protein genes as the reference, we picked the best hit
538 for each pigment gene in the three species, and then filtered those with less than 50% of matching
539 rate and less than 50% overall amino acid identity (along the entire sequence). Reciprocal
540 best-match pairs were defined as orthologs. Multiple sequence alignments (MSA) were
541 performed using MUSCLE.

542 Tyr (Tyrosinase) is a key rate-limiting enzyme in the melanin production. Mutations in this
543 gene cause human oculocutaneous albinism 1. By comparing *Tyr* ortholog sequences, we found a
544 nucleotide mutation (G419R) in one copy of Sa (also testified by PCR), which is identical with
545 one of human mutations in *Tyr* (<http://www.ifpcs.org/albinism/oca1mut.html>). This mutation site
546 has been first reported in Caucasian patients [72], then tested in Indo-Pakistani [73], Indian [74]
547 and Non-Hispanic Caucasians [75] patients. The alignment results are presented in [Figure S18](#).

548 We also checked melanin-associated-gene expression in the skin transcriptome of the three
549 *Sinocyclocheilus* species. Compared to Sg, *Tyr* (P-value of $1.69E-05$, FDR of $3.38E-05$) and
550 *Tyrp1* (P-value of $3.87E-67$, FDR of $1.55E-66$) in Sa was significantly down-regulated, while no

551 expression was detectable in Sr. Meanwhile, *Dct* expressed in Sr was significantly up-regulated
552 (Table S20). The expression levels of genes in the melanogenesis pathways in Sa are the lowest.

553

554 6.3 Scale development

555 Edar (Ectodysplasin-A receptor) is a member of the Tumor necrosis factor receptor superfamily
556 [76] which works together with other proteins during embryonic development and plays a key
557 role in initiation of hair development in mammals. It encoded a protein with approximately 450
558 amino acids, including front signal peptide, extracellular region, transmembrane region and
559 cytoplasmic region. Recently, Edar has been shown to be required for scale development in fishes
560 [77].

561 We downloaded *Edar* encoded protein (*ENSDARP00000045289*) from the zebrafish genome
562 and predicted two copies of *Edar* genes in the three *Sinocyclocheilus* genomes. Subsequently, we
563 aligned these proteins using MUSCLE against *ENSDARP00000045289* and found deletions
564 (Figure S20) in the *Edar* gene of Sa.

565 In our work, we found two copies of *Edar* genes in the genomes of the three *Sinocyclocheilus*
566 species while both of the two copies in Sa present deletions (also testified by PCR) including
567 signal peptide. These signal peptides guide new protein transfer across the membrane and parts of
568 the extracellular region. The deletions may make the functional protein incomplete resulting in
569 little scale covering on the skin surface of Sa (Figure S29).

570

571 6.4 Hearing

572 Using the same method as described in the pigment-related genes, we obtained a series of
573 genes associated with deafness and ear development. MSA was executed to find unique
574 amino-acid changes. *Mpv17* encodes a mitochondrial inner membrane protein associated with
575 mitochondrial DNA depletion syndrome. It has been reported that a *Mpv17*-deficient mice can
576 cause severe morphological degeneration of the cochlea, resulting in sensorineural hearing loss at
577 2-months-old [78]. We took out the best hit for *Mpv17* and found that one copy of this gene in Sa
578 has a deletion in the signal region which may lead to the functional decline (Figure S22). *Ush2a*
579 is another gene that altered significantly in Sa. Two amino-acid changes, R334S [79, 80] and
580 V382A [81] (also shown by PCR) were predicted (Figure S19). It has been proved that the
581 encoded protein is found in the basement membrane, and may be important in the development of
582 the inner ear and retina. In addition, mutations within this gene have been associated with Usher
583 syndrome type IIa and retinitis pigmentosa
584 (https://grenada.lumc.nl/LOVD2/Usher_montpellier/home.php?select_db=USH2A).

585 A frequent polymorphism in the translation start codon of Otoraplin (*Otor*) can abolish
586 translation and may be associated with forms of deafness. The encoded protein has
587 growth-inhibitory activity in melanoma cell lines [82]. Since we detected *Otor* in Sr, Sa and
588 *Astyanax mexicanus*, the *Otor* may have an impact on hearing and pigment development in the
589 cavefishes. Interestingly, the reconstructions of the saccular otolith morphology using
590 synchrotron X-ray microtomography (Figure 4), as well as both the anatomy of the swim bladder
591 (Figure S30) and distributions of the neuromasts of the trunk lateral line system (Figure S31),
592 indicated that these hearing related organs in Sa presented different degrees of weakness
593 compared to Sg and Sr. Furthermore, a comparison to the distribution of neuromasts on the head
594 (Figure 4) also hinted that the response to vibration of these three species was Sg>Sr>Sa.

595

596 6.5 Immune response

597 Protein sequences of 2,157 human immune response genes were downloaded from GO
598 database. A blsatp analysis was used with the e-value cutoff of 1e-5 to identify homolog genes in
599 zebrafish and the three *Sinocyclocheilus* species. Then, the following steps were used to filter
600 those candidate genes: a) deleted genes with identity and align rate less than 30%; b) executed
601 GO enrichment and removed genes that are not immune-related GO terms; c) discarded genes
602 that appeared with many immune-related GO terms. At last, we acquired the specific information
603 of immune genes of each *Sinocyclocheilus* species (Table S21).

604 The immune ability of cavefishes was suspected to be lower than surface fishes. The *Tlr* gene
605 family is a group of innate immune specific receptors. We analyzed the *Tlr* gene family among
606 the three *Sinocyclocheilus* species and zebrafish, and found that *Tlr* 5 in Sa is contracted (Table
607 S22). We performed similarity analysis with eleven sequenced fishes and representative species
608 of four branches (mammals: human and mouse, birds: chicken, amphibians: toad, reptile: anolis
609 and turtle) (Figure S27). The data showed that the fishes' *Tlr* copy numbers are greater than those
610 found in the other branches, and the *Sinocyclocheilus Tlr* genes are closer to those of zebrafish,
611 which is consistent with their close phylogenomic relationship.

612

613 6.6 Circadian rhythm

614 Previous research has shown that cavefishes lack circadian rhythms because they live in
615 perpetual darkness [83]. Mexican blind cavefish, *Astyanax mexicanus*, has a phenotype with
616 reduced sleep in comparison to surface populations [84]. We checked the rhythm pathway gene

617 copy number and their expressions in the semi-cave fishes Sr and cavefish Sa. Related zebrafish
618 genes were downloaded from the KEGG pathway database (map 04710) and aligned to the gene
619 sets of the three *Sinocyclocheilus* species. The homolog genes of these species were determined
620 according to the identity with zebrafish's category and their copy numbers were counted (Table
621 S27). Homolog protein sequence multiple alignments were performed with MUSCLE and the
622 structural variations were validated. The data demonstrated that both copies of Skp1 protein in Sa
623 were deleted in the N-terminal (Figure S24) (also testified by PCR). It was reported that Skp1
624 protein serves as an adapter in SCF complex, which links the F-box protein to Cull1 and plays an
625 important role in the maintenance of animal circadian rhythm. Meanwhile, transcriptomic
626 analysis showed that the rhythm genes' expression levels in the three *Sinocyclocheilus* species
627 declined in the order Sg>Sr>Sa (Figure S26). These observations together suggest a weak
628 circadian rhythm in the cavefish Sa.

629

630

631 6.7 Sense of taste

632 There are five primary tastes: sweet, bitter, salty, sour and umami. Peptides of Tas1r1 and
633 Tas1r3 form a heterodimer to detect the umami taste, whereas peptides of Tas1r2 and Tas1r3 form
634 a heterodimer for sweet detection. *Tas2r* genes are responsible for the bitter taste [85].
635 Furthermore, *ENaC* and *Pkd211*, channels responsible for sensing salty and sour tastes,
636 respectively. In the three *Sinocyclocheilus* genomes, the copy numbers of *Tas1r3* and *Pkd211* were
637 double those in the zebrafish genome. In particular, the copy number of *Tas1r2* (responsible for
638 sweet) was more than double compared with zebrafish, suggesting that the *Sinocyclocheilus*
639 fishes were more sensitive to sweet (Table S28). The *ENaC* gene was not identified in the three
640 *Sinocyclocheilus* genomes and this may suggest that *Sinocyclocheilus* fishes have lost the ability
641 to taste salt.

642 In conclusion, the genomic data suggest that the *Sinocyclocheilus* fishes are likely able to sense
643 umami, sweet, bitter, and sour tastes, but not salty, which is consistent with their current
644 freshwater environment.

645 **References (1-85)**

- 646 1. Zhao Y, Zhang C. Endemic Fishes of Sinocyclocheilus (Cypriniformes: Cyprinidae) in China-Species
647 diversity, Cave adaptation, Systematics and Zoogeography. Beijing: Science Press; 2009.
- 648 2. Li R, Fan W, Tian G, Zhu H, He L, Cai J et al. The sequence and de novo assembly of the giant panda
649 genome. *Nature*. 2010; 463(7279):311-7.
- 650 3. Li R, Yu C, Li Y, Lam T-W, Yiu S-M, Kristiansen K et al. SOAP2: An improved ultrafast tool for short
651 read alignment. *Bioinformatics*. 2009; 25(15):1966-7.
- 652 4. Wang Z, Pascual-Anaya J, Zadissa A, Li W, Niimura Y, Huang Z et al. The draft genomes of soft-shell
653 turtle and green sea turtle yield insights into the development and evolution of the turtle-specific body
654 plan. *Nat Genet*. 2013; 45(6):701-6.
- 655 5. Parra G, Bradnam K, Korf I. CEGMA: a pipeline to accurately annotate core genes in eukaryotic
656 genomes. *Bioinformatics*. 2007; 23(9):1061-7.
- 657 6. Parra G, Bradnam K, Ning Z, Keane T, Korf I. Assessing the gene space in draft genomes. *Nucleic Acids*
658 *Res*. 2009; 37(1):289-97.
- 659 7. Jiang Y, Xie M, Chen W, Talbot R, Maddox JF, Faraut T et al. The sheep genome illuminates biology of
660 the rumen and lipid metabolism. *Science*. 2014; 344(6188):1168-73.
- 661 8. Tarailo-Graovac M, Chen N. Using RepeatMasker to identify repetitive elements in genomic sequences.
662 *Curr Protoc Bioinformatics*. 2009; Chapter 4:Unit 4.10.
- 663 9. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J. Repbase Update, a
664 database of eukaryotic repetitive elements. *Cytogenet Genome Res*. 2005; 110(1-4):462-7.
- 665 10. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res*. 1999;
666 27(2):573-80.
- 667 11. Elsik CG, Mackey AJ, Reese JT, Milshina NV, Roos DS, Weinstock GM. Creating a honey bee
668 consensus gene set. *Genome Biol*. 2007; 8(1):R13.
- 669 12. Cho YS, Hu L, Hou H, Lee H, Xu J, Kwon S et al. The tiger genome and comparative analysis with lion
670 and snow leopard genomes. *Nat Commun*. 2013; 4:2433.
- 671 13. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol*.
672 1990; 215(3):403-10.
- 673 14. Birney E, Clamp M, Durbin R. GeneWise and Genomewise. *Genome Res*. 2004; 14(5):988-95.
- 674 15. Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B. AUGUSTUS: ab initio prediction of
675 alternative transcripts. *Nucleic Acids Res*. 2006; 34(Web Server issue):W435-9.
- 676 16. Korf I. Gene finding in novel genomes. *BMC Bioinformatics*. 2004; 5:59.
- 677 17. Majoros WH, Pertea M, Salzberg SL. TigrScan and GlimmerHMM: two open source ab initio
678 eukaryotic gene-finders. *Bioinformatics*. 2004; 20(16):2878-9.
- 679 18. Trapnell C, Pachter L, Salzberg SL. TopHat: discovering splice junctions with RNA-Seq.
680 *Bioinformatics*. 2009; 25(9):1105-11.
- 681 19. Trapnell C, Hendrickson DG, Sauvageau M, Goff L, Rinn JL, Pachter L. Differential analysis of gene

- 682 regulation at transcript resolution with RNA-seq. *Nat Biotechnol.* 2013; 31(1):46-53.
- 683 20. Attwood TK, Beck ME. PRINTS--a protein motif fingerprint database. *Protein Eng.* 1994; 7(7):841-8.
- 684 21. Hulo N, Bairoch A, Bulliard V, Cerutti L, De Castro E, Langendijk-Genevaux PS et al. The PROSITE
685 database. *Nucleic Acids Res.* 2006; 34(Database issue):D227-30.
- 686 22. Finn RD, Mistry J, Tate J, Coggill P, Heger A, Pollington JE et al. The Pfam protein families database.
687 *Nucleic Acids Res.* 2010; 38(Database issue):D211-22.
- 688 23. Bru C, Courcelle E, Carrere S, Beausse Y, Dalmar S, Kahn D. The ProDom database of protein domain
689 families: more emphasis on 3D. *Nucleic Acids Res.* 2005; 33(Database issue):D212-5.
- 690 24. Letunic I, Doerks T, Bork P. SMART 7: recent updates to the protein domain annotation resource.
691 *Nucleic Acids Res.* 2012; 40(Database issue):D302-5.
- 692 25. Mi H, Lazareva-Ulitsky B, Loo R, Kejariwal A, Vandergriff J, Rabkin S et al. The PANTHER database
693 of protein families, subfamilies, functions and pathways. *Nucleic Acids Res.* 2005; 33(Database
694 issue):D284-8.
- 695 26. Zdobnov EM, Apweiler R. InterProScan--an integration platform for the signature-recognition methods
696 in InterPro. *Bioinformatics.* 2001; 17(9):847-8.
- 697 27. Harris MA, Clark J, Ireland A, Lomax J, Ashburner M, Foulger R et al. The Gene Ontology (GO)
698 database and informatics resource. *Nucleic Acids Res.* 2004; 32(Database issue):D258-61.
- 699 28. Kanehisa M, Goto S, Shuichi K, Akihiro N. The KEGG databases at GenomeNet. *Nucleic Acids Res.*
700 2002; 30(1):42-6.
- 701 29. Griffiths-Jones S, Bateman A, Marshall M, Khanna A, Eddy SR. Rfam: an RNA family database.
702 *Nucleic Acids Res.* 2003; 31(1):439-41.
- 703 30. Gardner PP, Daub J, Tate JG, Nawrocki EP, Kolbe DL, Lindgreen S et al. Rfam: updates to the RNA
704 families database. *Nucleic Acids Res.* 2009; 37(Database issue):D136-40.
- 705 31. Lowe TM, Eddy SR. tRNAscan-SE: a program for improved detection of transfer RNA genes in
706 genomic sequence. *Nucleic Acids Res.* 1997; 25(5):955-64.
- 707 32. Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms and methods
708 to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. *Syst Biol.*
709 2010; 59(3):307-21.
- 710 33. Guindon S, Delsuc F, Dufayard JF, Gascuel O. Estimating maximum likelihood phylogenies with
711 PhyML. *Methods Mol Biol.* 2009; 537:113-37.
- 712 34. Huelsenbeck JP, Ronquist F. MRBAYES: Bayesian inference of phylogenetic trees. *Bioinformatics.*
713 2001; 17(8):754-5.
- 714 35. Zhang G, Cowled C, Shi Z, Huang Z, Bishop-Lilly KA, Fang X et al. Comparative analysis of bat
715 genomes provides insight into the evolution of flight and immunity. *Science.* 2013; 339(6118):456-60.
- 716 36. Donoghue PC, Benton MJ. Rocks and clocks: calibrating the Tree of Life using fossils and molecules.
717 *Trends Ecol Evol.* 2007; 22(8):424-31.
- 718 37. Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol.* 2007; 24(8):1586-91.

- 719 38. Yang Z. PAML: a program package for phylogenetic analysis by maximum likelihood. *Comput Appl*
720 *Biosci.* 1997; 13(5):555-6.
- 721 39. Wikstrom N, Savolainen V, Chase MW. Evolution of the angiosperms: calibrating the family tree. *Proc*
722 *Biol Sci.* 2001; 268(1482):2211-20.
- 723 40. Wang ZQ. A new Permian gnetalean cone as fossil evidence for supporting current molecular phylogeny.
724 *Ann Bot.* 2004; 94(2):281-8.
- 725 41. Li H, Durbin R. Inference of human population history from individual whole-genome sequences.
726 *Nature.* 2011; 475(7357):493-6.
- 727 42. Graur D, Li W-H. *Fundamentals of molecular evolution.* Sinauer Associates Sunderland; 2000.
- 728 43. Philipp JK. *Gnuplot in Action: Understanding Data with Graphs.* Manning Publications; 2009.
- 729 44. De Bie T, Cristianini N, Demuth JP, Hahn MW. CAFE: a computational tool for the study of gene
730 family evolution. *Bioinformatics.* 2006; 22(10):1269-71.
- 731 45. Hahn MW, De Bie T, Stajich JE, Nguyen C, Cristianini N. Estimating the tempo and mode of gene
732 family evolution from comparative genomic data. *Genome Res.* 2005; 15(8):1153-60.
- 733 46. Demuth JP, De Bie T, Stajich JE, Cristianini N, Hahn MW. The evolution of mammalian gene families.
734 *PLoS One.* 2006; 1:e85.
- 735 47. Yamashita T, Ono K, Ohuchi H, Yumoto A, Gotoh H, Tomonari S et al. Evolution of mammalian Opn5
736 as a specialized UV-absorbing pigment by a single amino acid mutation. *J Biol Chem.* 2014;
737 289(7):3991-4000.
- 738 48. Sun YB, Zhou WP, Liu HQ, Irwin DM, Shen YY, Zhang YP. Genome-wide scans for candidate genes
739 involved in the aquatic adaptation of dolphins. *Genome Biol Evol.* 2013; 5(1):130-9.
- 740 49. Zhang J, Nielsen R, Yang Z. Evaluation of an improved branch-site likelihood method for detecting
741 positive selection at the molecular level. *Mol Biol Evol.* 2005; 22(12):2472-9.
- 742 50. Huang da W, Sherman BT, Lempicki RA. Systematic and integrative analysis of large gene lists using
743 DAVID bioinformatics resources. *Nat Protoc.* 2009; 4(1):44-57.
- 744 51. Abdel-Salam GM, Hassan NA, Kayed HF, Aligianis IA. Phenotypic variability in Micro syndrome:
745 report of new cases. *Genet Couns.* 2007; 18(4):423-35.
- 746 52. Dursun F, Guven A, Morris-Rosendahl D. Warburg Micro syndrome. *J Pediatr Endocrinol Metab.* 2012;
747 25(3-4):379-82.
- 748 53. Riise R, Tornqvist K, Wright AF, Mykytyn K, Sheffield VC. The phenotype in Norwegian patients with
749 Bardet-Biedl syndrome with mutations in the BBS4 gene. *Arch Ophthalmol.* 2002; 120(10):1364-7.
- 750 54. Ohtsubo M, Sato M, Hikoya A, Hosono K, Minoshima S, Hotta Y. Case of Japanese patient with
751 x-linked ocular albinism associated with GPR143 gene mutation. *Jpn J Ophthalmol.* 2010;
752 54(6):624-6.
- 753 55. Bang YL, Nguyen TT, Trinh TT, Kim YJ, Song J, Song YH. Functional analysis of mutations in
754 UDP-galactose-4-epimerase (GALE) associated with galactosemia in Korean patients using
755 mammalian GALE-null cells. *FEBS J.* 2009; 276(7):1952-61.

- 756 56. Sandrock K, Bartsch I, Rombach N, Schmidt K, Nakamura L, Hainmann I et al. Compound
757 heterozygous mutations in 2 siblings with Hermansky-Pudlak syndrome type 1 (HPS1). *Klin Padiatr.*
758 2010; 222(3):168-74.
- 759 57. Charizopoulou N, Lelli A, Schraders M, Ray K, Hildebrand MS, Ramesh A et al. Gipc3 mutations
760 associated with audiogenic seizures and sensorineural hearing loss in mouse and human. *Nat Commun.*
761 2011; 2:201.
- 762 58. Meyer A, Van de Peer Y. From 2R to 3R: evidence for a fish-specific genome duplication (FSGD).
763 *Bioessays.* 2005; 27(9):937-45.
- 764 59. Adham IM, Eck TJ, Mierau K, Muller N, Sallam MA, Paprotta I et al. Reduction of spermatogenesis
765 but not fertility in Creb3l4-deficient mice. *Mol Cell Biol.* 2005; 25(17):7657-64.
- 766 60. Casalena G, Krick S, Daehn I, Yu L, Ju W, Shi S et al. Mpv17 in mitochondria protects podocytes
767 against mitochondrial dysfunction and apoptosis in vivo and in vitro. *Am J Physiol Renal Physiol.*
768 2014; 306(11):F1372-80.
- 769 61. Chang CY, Lin SC, Su WH, Ho CM, Jou YS. Somatic LMCD1 mutations promoted cell migration and
770 tumor metastasis in hepatocellular carcinoma. *Oncogene.* 2012; 31(21):2640-52.
- 771 62. Bagley RG, Rouleau C, Weber W, Mehraein K, Smale R, Curiel M et al. Tumor endothelial marker 7
772 (TEM-7): a novel target for antiangiogenic therapy. *Microvasc Res.* 2011; 82(3):253-62.
- 773 63. Murshed F, Farhana L, Dawson MI, Fontana JA. NF-kappaB p65 recruited SHP regulates
774 PDCD5-mediated apoptosis in cancer cells. *Apoptosis.* 2014; 19(3):506-17.
- 775 64. Green MR, Vicente-Duenas C, Romero-Camarero I, Long Liu C, Dai B, Gonzalez-Herrero I et al.
776 Transient expression of Bcl6 is sufficient for oncogenic function and induction of mature B-cell
777 lymphoma. *Nat Commun.* 2014; 5:3904.
- 778 65. Kuijpers MJ, van der Meijden PE, Feijge MA, Mattheij NJ, May F, Govers-Riemslog J et al. Factor XII
779 Regulates the Pathological Process of Thrombus Formation on Ruptured Plaques. *Arterioscler Thromb*
780 *Vasc Biol.* 2014; 34(8):1674-80.
- 781 66. Chuang SS, Helvig C, Taimi M, Ramshaw HA, Collop AH, Amad M et al. CYP2U1, a novel human
782 thymus- and brain-specific cytochrome P450, catalyzes omega- and (omega-1)-hydroxylation of fatty
783 acids. *J Biol Chem.* 2004; 279(8):6305-14.
- 784 67. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic*
785 *Acids Res.* 2004; 32(5):1792-7.
- 786 68. Meng F, Braasch I, Phillips JB, Lin X, Titus T, Zhang C et al. Evolution of the eye transcriptome under
787 constant darkness in *Sinocyclocheilus cavefish*. *Mol Biol Evol.* 2013; 30(7):1527-43.
- 788 69. Blanton SH, Cortez A, Stal S, Mulliken JB, Finnell RH, Hecht JT. Variation in IRF6 contributes to
789 nonsyndromic cleft lip and palate. *Am J Med Genet A.* 2005; 137A(3):259-62.
- 790 70. Bangal S, Bhandari A, Dhayதாக P, Gogri P. Gyrate atrophy of choroid and retina with myopia,
791 cataract and systemic proximal myopathy: A rare case report from rural India. *Australas Med J.* 2012;
792 5(12):639-42.

- 793 71. Pittler SJ, Baehr W, Wasmuth JJ, McConnell DG, Champagne MS, vanTuinen P et al. Molecular
794 characterization of human and bovine rod photoreceptor cGMP phosphodiesterase alpha-subunit and
795 chromosomal localization of the human gene. *Genomics*. 1990; 6(2):272-83.
- 796 72. King RA, Mentink MM, Oetting WS. Non-random distribution of missense mutations within the human
797 tyrosinase gene in type I (tyrosinase-related) oculocutaneous albinism. *Mol Biol Med*. 1991;
798 8(1):19-29.
- 799 73. Tripathi RK, Bunday S, Musarella MA, Droetto S, Strunk KM, Holmes SA et al. Mutations of the
800 tyrosinase gene in Indo-Pakistani patients with type I (tyrosinase-deficient) oculocutaneous albinism
801 (OCA). *Am J Hum Genet*. 1993; 53(6):1173-9.
- 802 74. Chaki M, Sengupta M, Mukhopadhyay A, Subba Rao I, Majumder PP, Das M et al. OCA1 in different
803 ethnic groups of india is primarily due to founder mutations in the tyrosinase gene. *Ann Hum Genet*.
804 2006; 70(Pt 5):623-30.
- 805 75. Hutton SM, Spritz RA. Comprehensive analysis of oculocutaneous albinism among non-Hispanic
806 caucasians shows that OCA1 is the most prevalent OCA type. *J Invest Dermatol*. 2008;
807 128(10):2442-50.
- 808 76. Monreal AW, Ferguson BM, Headon DJ, Street SL, Overbeek PA, Zonana J. Mutations in the human
809 homologue of mouse *dl* cause autosomal recessive and dominant hypohidrotic ectodermal dysplasia.
810 *Nat Genet*. 1999; 22(4):366-9.
- 811 77. Kondo S, Kuwahara Y, Kondo M, Naruse K, Mitani H, Wakamatsu Y et al. The medaka *rs-3* locus
812 required for scale development encodes ectodysplasin-A receptor. *Curr Biol*. 2001; 11(15):1202-6.
- 813 78. Muller M, Smolders JW, Meyer zum Gottesberge AM, Reuter A, Zwacka RM, Weiher H et al. Loss of
814 auditory function in transgenic *Mpv17*-deficient mice. *Hear Res*. 1997; 114(1-2):259-63.
- 815 79. Dreyer B, Tranebjaerg L, Rosenberg T, Weston MD, Kimberling WJ, Nilssen O. Identification of novel
816 *USH2A* mutations: implications for the structure of *USH2A* protein. *Eur J Hum Genet*. 2000;
817 8(7):500-6.
- 818 80. Baux D, Larriue L, Blanchet C, Hamel C, Ben Salah S, Vielle A et al. Molecular and in silico analyses
819 of the full-length isoform of usherin identify new pathogenic alleles in Usher type II patients. *Hum*
820 *Mutat*. 2007; 28(8):781-9.
- 821 81. Garcia-Garcia G, Aparisi MJ, Jaijo T, Rodrigo R, Leon AM, Avila-Fernandez A et al. Mutational
822 screening of the *USH2A* gene in Spanish *USH* patients reveals 23 novel pathogenic mutations.
823 *Orphanet J Rare Dis*. 2011; 6:65.
- 824 82. Robertson NG, Heller S, Lin JS, Resendes BL, Weremowicz S, Denis CS et al. A novel conserved
825 cochlear gene, *OTOR*: identification, expression analysis, and chromosomal mapping. *Genomics*.
826 2000; 66(3):242-8.
- 827 83. Cavallari N, Frigato E, Vallone D, Frohlich N, Lopez-Olmeda JF, Foa A et al. A blind circadian clock in
828 cavefish reveals that opsins mediate peripheral clock photoreception. *PLoS Biol*. 2011; 9(9):e1001142.
- 829 84. Duboue ER, Keene AC, Borowsky RL. Evolutionary convergence on sleep loss in cavefish populations.

- 830 Curr Biol. 2011; 21(8):671-6.
831 85. Yarmolinsky DA, Zuker CS, Ryba NJ. Common sense about taste: from mammals to insects. Cell. 2009;
832 139(2):234-44.