

1 **Chromosome-level genome assembly of the African pike,**

2 *Hepsetus odoe*

3

4 Xiao Du^{1,2,3,*}, Xiaoning Hong^{1,2,3,4,*}, Guangyi Fan^{1,2,3,*}, Xiaoyun Huang^{1,2,3}, Shuai
5 Sun^{1,2,3}, Ouyang Bingjie^{1,2,3}, He Zhang^{1,2,3}, Mengqi Zhang^{1,2,3}, Shanshan Liu^{1,2,3}, Xin
6 Liu^{1,2,3,#} & Wenwei Zhang^{2,#}

7

8 ¹ BGI-Qingdao, BGI-Shenzhen, Qingdao, 266555, China

9 ² BGI-Shenzhen, Shenzhen, 518083, China

10 ³ China National GeneBank, BGI-Shenzhen, Shenzhen, 518120, China

11 ⁴ BGI Education Center, University of Chinese Academy of Sciences, Shenzhen, 236009,
12 China

13 * These authors contributed equally: Xiao Du, Xiaoning Hong, and Guangyi Fan

14

15

16 # Correspondence authors: Xin Liu (liuxin@genomics.cn); Wenwei Zhang
17 (zhangww@genomics.cn)

18 **Abstract**

19 The order Characiformes is one of the largest components of the freshwater teleost fauna
20 inhabiting exclusively in South America and Africa with great ecological and economical
21 significance. Yet, quite limited genomic resources are available to study this group and
22 their transatlantic vicariance. In this study we present a chromosome-level genome
23 assembly of the African pike (*Hepsetus odoe*), a representative member of the African
24 Characiformes. To this end, we generated 119, 11, and 67 Gb reads using the single tube
25 long fragment read (stLFR), Oxford Nanopore, and Hi-C sequencing technologies,
26 respectively. We obtained an 862.1 Mb genome assembly with the contig and scaffold
27 N50 of 347.4 kb and 25.8 Mb, respectively. Hi-C sequencing produced 29 chromosomes
28 with 742.5 Mb, representing 86.1% of the genome. 24,314 protein-coding genes were
29 predicted and 23,999 (98.7%) genes were functionally annotated. The chromosomal-scale
30 genome assembly will be useful for functional and evolutionary studies of the African
31 pike and promote the study of Characiformes speciation and evolution.

32

33

34 **Background & Summary**

35 The order Characiformes is one of the largest components of the freshwater fish fauna
36 worldwide, comprising about 2,000 ecologically and morphologically diverse fish living
37 in rivers and lakes exclusively in Africa and South America. Characiformes are
38 ecologically important by playing crucial roles in energy flux and material cycling in
39 river systems¹. Moreover, they have great significance for local economy because of diet
40 component of livestock and humans². As Characiformes are exclusively freshwater fishes,
41 their transatlantic distribution was proposed ascribed to the split of South America and
42 Africa in the Early Cretaceous fragmentation of western Gondwana³. This distribution
43 across the Atlantic Ocean displays asymmetry in the number of species, with *circa* 220
44 reported species in Africa and over 1,700 species in South America. High fragmentation
45 has been reported in the South American species, compared to the lower fragmentation
46 and variability in African ones^{4,5}. However, quite limited genomic information has been
47 available to study the Characiformes vicariance. Despite the large amount and high
48 diversity of species, presently only three genomes from two families *Characidae*

49 (*Astyanax mexicanus*) and *Serrasalminae* (*Pygocentrus nattereri*, *Colossoma*
50 *macropomum*) have been released in Characiformes, which all belong to the South
51 American lineages. No genomes of African Characiformes have been reported. Genomic
52 studies of African Characiformes would highly promote the understanding of
53 Characiformes evolution and speciation during the continent fragmentation.
54
55 African pike, *Hepsetus odoe*, is a representative African Characiformes that belongs to
56 the family *Hepsetidae*. It is a torpedo-shaped predatory and piscivorous species
57 distributed in the freshwater basins in central and western Africa⁶, bearing a striking
58 resemblance to the European pike. One of the most striking features is their dentition
59 with the lower jaw filled with two rows of sharp point teeth while the upper with only
60 one row. Due to roles in freshwater food chain and diet component of livestock and
61 humans, *H. odoe* are biologically and economically important². *H. odoe* was reported the
62 only species in the *Hepsetidae* family, until recently five additional species were
63 described by recent studies^{6,7}. Despite of the high economic and evolutionary importance,
64 no genome data are available for this group.
65
66 A high-quality genome assembly of the African pike will highly facilitate the study of its
67 functional and evolutionary genomics, which also will promote the understanding of other
68 African Characiformes along with their divergence from the South American
69 Characiformes. Therefore, in this study we report a chromosome-scale genome assembly
70 of *H. odoe* using single tube long fragment read (stLFR)⁸, Oxford Nanopore, and Hi-C
71 technologies (Additional file 1: Fig. S1). We obtained a genome assembly of 862.1 Mb
72 with the contig and scaffold N50 of 347.4 kb and 25.8 Mb, respectively. With
73 chromosome-level scaffolding, 29 scaffolds were constructed corresponding to 29
74 chromosomes with a total length of 742.5 Mb, representing 86.1% of all genome
75 sequences. 24,314 protein-coding genes were predicted in the assembly, and 98.7% of
76 them were functionally annotated. The chromosomal-level genome assembled here will
77 be useful for functional and evolutionary research of the African pike. It is the first
78 genome assembly in the African Characiformes and will promote the understanding of
79 Characiformes speciation and evolution.

80

81 **Methods**

82 **Sampling and sequencing**

83 Long genomic DNA (gDNA) from muscle tissue of a male African pike was isolated
84 using a conventional approach for sufficient DNA quality⁹. DNA integrity was checked
85 using agarose gel electrophoresis. The sequencing libraries were constructed via stLFR
86 technology according to the standard protocol via the MGIEasy stLFR library preparation
87 kit (PN:1000005622)¹⁰ and were sequenced on BGISEQ-500 platform. To overcome the
88 gaps (long ambiguous sequences) induced by repeats, library preparation and sequencing
89 were performed on the MinION nanopore sequencer (Oxford Nanopore Technologies,
90 Oxford, UK) for generating long reads, following the base protocols of Oxford Nanopore.
91 To get a high-resolution genome contact map, we used *in situ* Hi-C according to the
92 protocol of previous study with some modifications¹¹. The restriction endonuclease MboI
93 was used to digest DNA, followed by biotinylated residue labeling. The Hi-C library was
94 sequenced on BGISEQ-500 platform with 100 bp pair-end sequencing.

95

96 **Ethics statement**

97 The adult male African pike was purchased from the fish and aquarium market in
98 Guangzhou, Guangdong Province, China in May 2018. The experimental procedures
99 followed the guidelines approved by the institutional review board on bioethics and
100 biosafety of BGI (IRB-BGI). The experiment was authorized by the IRB-BGI (under NO.
101 FT17007). The review procedures in IRB-BGI meet good clinical practice (GCP)
102 principles.

103

104 **De novo assembly, and chromosome construction**

105 The k-mer frequency distribution analysis¹² was used to estimate the African pike
106 genome size. According to the 17-mer analysis, the genome size of African pike was
107 estimated to be 995 Mb (Table 2; Additional file 1: Fig. S2).

108

109 We obtained 118.6 Gb (~141X; Table 1) raw sequencing reads from stLFR. We used
110 SOAPfilter v.2.2, a package in SOAPdenovo2¹³ to filter reads with low quality reads (>

111 40% low-quality bases, $Q < 7$), PCR duplication, or adapter contamination. After filtering,
112 60.4 Gb (~72X; Table 1) clean reads were obtained for genome assembly. Supernova
113 assembler v2.0.1 (10X Genomics, Pleasanton, CA) was used to build contigs and
114 scaffolds, and gaps were closed by GapCloser (v1.2)¹³. With stLFR data, the generated
115 African pike assembly was 859.2 Mb. The contig and scaffold N50 were 43.9 kb and 5.1
116 Mb, respectively (Table 3). On basis of that, we generated a total of 11.0 Gb (~13X;
117 Table 1) long reads on the MinION nanopore sequencer to further fill the gaps using
118 TGSGapFiller¹⁴ with default parameters. After gap filling, the contig length was highly
119 elevated with contig N50 of 352.1 kb (Table 3).

120

121 Reads from Hi-C library¹⁵ were used to generate a chromosomal-level genome assembly.
122 First, we obtained 65.8 Gb (~76X, Table 1) clean sequencing data from the Hi-C library
123 by removing reads containing more than 1% unidentified (N) bases and low-quality bases
124 (quality value < 10) using SOAPnuke (v1.5.4)¹⁶ with parameters “-l 10 -q 0.1 -n 0.01 -Q
125 2”. Next, we used HiC-Pro pipeline (v2.8.0)¹⁷ for quality control to generate valid reads.
126 Of all 658,260,000 raw pair-end reads, there were 22.78% valid (149,912,370) paired Hi-
127 C reads suitable for following analysis. We used Juicer (v.1.5)¹⁸, an open-source and
128 fully-automated pipeline for pretreatment of Hi-C datasets, for analyzing valid Hi-C
129 datasets and producing the alignment result. Lastly, we applied 3D-DNA workflow (3D
130 *de novo* assembly, v.170123)¹⁹ to create the ordered-and-oriented genome sequences in
131 chromosome level with the main parameter “-m haploid -s 4 -c 29”. We assembled 29
132 chromosomes of *H. odoe* ranging from 8.03 Mb to 34.85 Mb with the total length of
133 742.5 Mb (Table 4; Fig. 1), which possessed 86.1% of all genome sequences. The final
134 African pike genome assembly spanned 862.1 Mb and 29 chromosomes, accounting for
135 86.6 % of the estimated genome size, with contig and scaffold N50 of 347.4 kb and 25.8
136 Mb, respectively. The constructed 29 chromosomes agreed with the previous karyotype
137 analysis of *H. odoe*⁶.

138

139 **Gene prediction and functional annotation**

140 To facilitate gene prediction in the genome, repetitive elements were identified first. Two
141 methods (*de novo* and homology-based predictions) were performed in the repeat

142 annotation of the African pike genome. In the *de novo* method, a *de novo* library was
143 built via running RepeatModeler (v1.0.8)²⁰ and LTR-FINDER (v1.0.6)²¹, and the
144 predicted model was applied to identify interspersed repetitive elements by
145 RepeatMasker (v4.0.5). In the homology-based prediction, detection of interspersed
146 repeats was realized by aligning the genome against the Repbase database²² at DNA and
147 protein levels using RepeatMasker and RepeatProteinMask (v4.0.5)²³. Tandem repeats
148 were predicted by TRF (v4.07). By integrating results of above approaches, 317.3 Mb
149 repetitive sequences were predicted, representing 36.7% of the genome assembly (Table
150 5). Finally, 284.1 Mb TEs were identified, accounting for 32.9% of the genome assembly.
151 The repetitive element annotations were summarized in Table 6. Those repetitive
152 sequences were masked to reduce the interference for the following gene predictions.
153
154 Next, we conducted structural and functional annotation for the assembled genome. For
155 structural annotation, both homology-based and *de novo* prediction approaches were
156 applied. In *de novo* prediction, AUGUSTUS (v3.1)²⁴ and GENSCAN (v2009)²⁵ were
157 utilized to predict the gene model with zebrafish data as a training set, and 23,163 and
158 29,084 protein-coding genes were predicted, respectively (Table 7). The homology-based
159 prediction of genome assembly was realized by referring to the NCBI protein repertoires
160 of six homologous species including Mexican tetra (*Astyanax mexicanus*), red-bellied
161 piranha (*Pygocentrus nattereri*), channel catfish (*Ictalurus punctatus*), common carp
162 (*Cyprinus carpio*), iridescent shark (*Pangasianodon hypophthalmus*), and zebrafish
163 (*Danio rerio*). After mapping the protein sequences to the repeat-masked African pike
164 genome using BLAST²⁶ (*E*-value cutoff of 1×10^{-5}), GeneWise (v2.4.1)²⁷ was used to
165 predict gene models by aligning homologous genome sequences against the matched
166 proteins. Lastly, we performed GLEAN to integrate all above gene models and obtained a
167 non-redundant gene set consisting of 24,314 protein-coding genes (Table 7). There were
168 9.77 exons per gene and the average length of coding sequences (CDS) was 1,712 bp
169 (Table 7). Gene function was annotated with TrEMBL²⁸, Swissprot²⁸, InterPro²⁹, Gene
170 Ontology³⁰, and Kyoto Encyclopedia of Genes and Genomes (KEGG)³¹ databases.
171 Ultimately, 23,999 genes (98.7% of the total) in African pike were functionally annotated
172 (Table 8).

173

174 **Genome features of the African pike**

175 CpG islands (CGIs), which are a significant group of CpG dinucleotide repeats in
176 genome regions, are functionally important for genomic studies. The CGIs were
177 identified across the genome using CpGIScan³². Ultimately, 24,297 CGIs were identified
178 with a total length of 15.5 Mb. A range of genome features including gene density, repeat
179 content, GC content, and GGI content were summarized and depicted in Fig. 2a. The
180 CpG density was found positively correlated with GC content, gene density, and repeat
181 content (Fig. 2b), following a similar pattern observed in other published fish and
182 mammals genomes³³⁻³⁵.

183

184 **Gene family identification**

185 Gene family analysis among species provides significant insights into phylogenetic and
186 evolutionary studies. The protein-coding genes from the African pike assembly, two
187 sequenced species in order Characiformes (*A. mexicanus* and *P. nattereri*), and five other
188 sequenced species including Atlantic salmon (*Salmo salar*), yellow catfish (*Tachysurus*
189 *fulvidraco*), northern pike (*Esox lucius*), electric eel (*Electrophorus electricus*), and
190 zebrafish (*Danio rerio*) were downloaded from NCBI database and analyzed. All-versus-
191 all protein similarities were computed using BLASTP²⁶ and the alignment results were
192 used by TreeFam (v4.0)³⁶ to deduce homologous gene sequences and identify gene
193 families. Orthologue clustering analysis of predicted genes was conducted using MCL
194 algorithm (Fig. 3a). Finally, we identified 9,661 gene families in the African pike genome
195 (Additional file 1: Table S1). Compared to the two South American Characiformes (*A.*
196 *mexicanus* and *P. nattereri*) and the zebrafish (*D. rerio*), 221 gene families were unique
197 in African pike (Fig. 3b).

198

199 **Phylogenetic analysis**

200 To study the evolutionary position of African pike, 3,106 single-copy genes from the
201 above seven species were used for constructing phylogenetic tree and estimating
202 divergence time. Protein sequences of single-copy gene families were aligned by
203 MUSCLE (v3.8.31)³⁷ and then were concatenated into a supergene matrix for each

204 species. The alignment results were processed into PhyML (v 3.0)^{37,38} to construct a ML
205 phylogenetic tree. Divergence time was inferred using the MCMCTree from the PAML
206 package³⁹. Divergence times from TimeTree database⁴⁰ were applied for calibration,
207 which include splits between *E. lucius* and *D. rerio* (198-211 Mya), between *D. rerio* and
208 *T. fulvidraco* (170-183 Mya), and between *T. fulvidraco* and *E. electricus* (122-136 Mya).
209 The phylogenetic tree showed that *H. odoe* was most closely related to *P. nattereri* with a
210 divergence time around 73.8 Mya, and together the clade formed the sister group to *A.*
211 *mexicanus* (Fig. 4a). The African family (*H. odoe*) was dispersed among the South
212 American families (*P. nattereri* than *A. mexicanus*). Although sharing similar pike-like
213 forms, the African pike was distantly related to the northern pike (*E. lucius*) with a
214 divergence time around 205 Mya.

215

216

217 **Expansion and contraction of gene families**

218 Based on the gene family clustering results and divergence time estimation, we used Café
219 (v2.1)⁴¹ to estimate the gene family expansion and contraction events during speciation.
220 Results showed that in African pike genome 769 gene families were found expanded and
221 1,346 gene families were contracted (Fig. 4a). The 284 significantly expanded and 81
222 significantly contracted gene families ($p < 0.05$; Additional file 1: Table S2) in African
223 pike were annotated with KEGG ortholog functions. Among that 93 (32.7%) gene
224 families were strikingly expanded in immune system (Fig. 4b).

225

226 **Data Records**

227 The sequencing data and genome assembly of the African pike were deposited in NCBI
228 under BioProject accession PRJNA625402. The datasets reported in this study are also
229 available in the CNGB Nucleotide Sequence Archive (CNSA: <https://db.cngb.org/cnsa>;
230 accession number CNP0001012).

231

232 **Technical Validation**

233 **Assessment of genome assembly**

234 The contig and scaffold N50 of the African pike genome were 347.4 kb and 25.8 Mb
235 respectively, with the longest scaffold 34,852,849 bp. We assessed the quality of the
236 assembled genome using the Benchmarking Universal Single-copy Orthologs (BUSCO
237 v3.0.2)⁴². The assembly reached 90.7% ~92.4% completeness compared to single-copy
238 ortholog gene sets from atinopterygii, metazoans, and vertebrates in BUSCOs (Table 9).
239 This demonstrates the high completeness of our genome assembly.

240

241 Next, we evaluated the assembly of the twenty-nine chromosomes. The genome assembly
242 was divided into 100kb bins. The signal for the interaction between any two bins was
243 defined by the count of Hi-C reads covered by those bins, and the signal intensities were
244 depicted in a heat map. The Hi-C heat map clearly split the bins into 29 blocks, and bins
245 within the same chromosome showed substantially larger signal intensities than bins
246 distributed on different chromosomes (Fig. 1). This demonstrates the high quality of the
247 chromosome assembly.

248

249 **Gene prediction and annotation validation**

250 Repetitive sequences in the assembly were masked before gene annotation. Gene model
251 prediction in the African pike was realized by using a combination of *de novo* and
252 homology-based approaches. Then the gene prediction results were integrated into a
253 consensus gene set by GLEAN. Annotation completeness of the African pike gene set
254 was assessed by BUSCO, reaching 92.5%~96.4% completeness (Table 9). In addition,
255 functional annotation of the predicted genes showed that 98.7% of them could be
256 assigned into at least one functional term (Table 8). These results clearly indicate the
257 annotated gene set is quite complete.

258

259 **Code Availability**

260 All commands used in the analysis were executed by following the manual of the
261 corresponding bioinformatics tools. There were no any custom specific codes.

262

263 **Acknowledgements**

264 This work was supported by the special funding of “Blue granary” scientific and
265 technological innovation of China (2018YFD0900301-05). We also thank for the
266 technical supports from China National Genebank in stLFR library construction and
267 sequencing.

268

269 **Author Contributions**

270 G. F., H. Z., and X. L. designed the study. G. F., X. L., and W. Z. supervised the study.
271 M. Z. and S. L. contributed to sample collection and sequencing experiments. X. D., X.
272 Hong, S. S., X. Huang and B.O. performed bioinformatics analyses. X. D., X. Hong, and
273 G. F. wrote the manuscript.

274

275 **Competing interests**

276 The authors declare no competing interests.

277 **References**

- 278 1 Taylor, B. W., Flecker, A. S. & Hall, R. O., Jr. Loss of a harvested fish species
279 disrupts carbon flow in a diverse tropical river. *Science* **313**, 833-836,
280 doi:10.1126/science.1128223 (2006).
- 281 2 Ogunola, O. S., Onada, O. A. & Falaye, A. E. Preliminary evaluation of some
282 aspects of the ecology (growth pattern, condition factor and reproductive biology)
283 of African pike, *Hepsetus odoe* (Bloch 1794), in Lake Eleiyele, Ibadan, Nigeria.
284 *Fisheries & Aquatic Science* **21**, 12.
- 285 3 Calcagnotto, D., Schaefer, S. A. & DeSalle, R. Relationships among characiform
286 fishes inferred from analysis of nuclear and mitochondrial gene sequences.
287 *Molecular phylogenetics and evolution* **36**, 135-153,
288 doi:10.1016/j.ympev.2005.01.004 (2005).
- 289 4 Orti, G. & Meyer, A. The radiation of characiform fishes and the limits of
290 resolution of mitochondrial ribosomal DNA sequences. *Syst Biol* **46**, 75-100,
291 doi:10.1093/sysbio/46.1.75 (1997).
- 292 5 Carvalho, P. C. *et al.* First Chromosomal Analysis in Hepsetidae (Actinopterygii,
293 Characiformes): Insights into Relationship between African and Neotropical Fish
294 Groups. *Frontiers in genetics* **8**, doi:10.3389/fgene.2017.00203 (2017).
- 295 6 Carvalho, P. C. *et al.* First Chromosomal Analysis in Hepsetidae (Actinopterygii,
296 Characiformes): Insights into Relationship between African and Neotropical Fish
297 Groups. *Frontiers in genetics* **8**, 203, doi:10.3389/fgene.2017.00203 (2017).
- 298 7 Decru, E., Vreven, E. & Snoeks, J. A revision of the Lower Guinean *Hepsetus*
299 species (Characiformes; Hepsetidae) with the description of *Hepsetus kingsleyae*
300 sp. nov. *Journal of fish biology* **82**, 1351-1375, doi:10.1111/jfb.12079 (2013).
- 301 8 Wang, O. *et al.* Efficient and unique cobarcoding of second-generation
302 sequencing reads from long DNA molecules enabling cost-effective and accurate
303 sequencing, haplotyping, and de novo assembly. *Genome Research* **29**, 798-808,
304 doi:10.1101/gr.245126.118 (2019).
- 305 9 Panova, M. *et al.* in *Marine genomics* (ed S. Bourlat) 13-44 (Humana Press,
306 2016).

- 307 10 Wang, O. *et al.* Efficient and unique cobarcoding of second-generation
308 sequencing reads from long DNA molecules enabling cost-effective and accurate
309 sequencing, haplotyping, and de novo assembly. *Genome Res* **29**, 798-808,
310 doi:10.1101/gr.245126.118 (2019).
- 311 11 Rao, S. S. *et al.* A 3D map of the human genome at kilobase resolution reveals
312 principles of chromatin looping. *Cell* **159**, 1665-1680,
313 doi:10.1016/j.cell.2014.11.021 (2014).
- 314 12 Li, R. *et al.* The sequence and de novo assembly of the giant panda genome.
315 *Nature* **463**, 311-317, doi:10.1038/nature08696 (2010).
- 316 13 Luo, R. *et al.* SOAPdenovo2: an empirically improved memory-efficient short-
317 read de novo assembler. *Gigascience* **1**, 18, doi:10.1186/2047-217X-1-18 (2012).
- 318 14 Guo, L. & Deng, L. 10.5281/zenodo.3446281 (2019).
- 319 15 Burton, J. N. *et al.* Chromosome-scale scaffolding of de novo genome assemblies
320 based on chromatin interactions. *Nature Biotechnology* **31**, 1119-1125,
321 doi:10.1038/nbt.2727 (2013).
- 322 16 Chen, Y. *et al.* SOAPnuke: a MapReduce acceleration-supported software for
323 integrated quality control and preprocessing of high-throughput sequencing data.
324 *Gigascience* **7**, gix120 (2017).
- 325 17 Servant, N. *et al.* HiC-Pro: an optimized and flexible pipeline for Hi-C data
326 processing. *Genome Biology* **16**, 259, doi:10.1186/s13059-015-0831-x (2015).
- 327 18 Durand, N. C. *et al.* Juicer Provides a One-Click System for Analyzing Loop-
328 Resolution Hi-C Experiments. *Cell systems* **3**, 95-98,
329 doi:10.1016/j.cels.2016.07.002 (2016).
- 330 19 Dudchenko, O. *et al.* De novo assembly of the *Aedes aegypti* genome using Hi-C
331 yields chromosome-length scaffolds. *Science* **356**, 92-95,
332 doi:10.1126/science.aal3327 (2017).
- 333 20 Smith, A. & Hubley, R. (2008-2015).
- 334 21 Xu, Z. & Wang, H. LTR_FINDER: an efficient tool for the prediction of full-
335 length LTR retrotransposons. *Nucleic Acids Res.* **35**, W265-W268,
336 doi:10.1093/nar/gkm286 (2007).

- 337 22 Bao, W., Kojima, K. K. & Kohany, O. Repbase Update, a database of repetitive
338 elements in eukaryotic genomes. *Mobile DNA* **6**, 11, doi:10.1186/s13100-015-
339 0041-9 (2015).
- 340 23 Tarailo-Graovac, M. & Chen, N. Using RepeatMasker to identify repetitive
341 elements in genomic sequences. *Current protocols in bioinformatics* **25**, 4-10
342 (2009).
- 343 24 Stanke, M. *et al.* AUGUSTUS: *ab initio* prediction of alternative transcripts.
344 *Nucleic Acids Research* **34**, W435-W439, doi:10.1093/nar/gkl200 (2006).
- 345 25 Burge, C. & Karlin, S. Prediction of complete gene structures in human genomic
346 DNA. *Journal of Molecular Biology* **268**, 78-94, doi:10.1006/jmbi.1997.0951
347 (1997).
- 348 26 Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local
349 alignment search tool. *Journal of Molecular Biology* **215**, 403-410,
350 doi:10.1016/S0022-2836(05)80360-2 (1990).
- 351 27 Birney, E., Clamp, M. & Durbin, R. GeneWise and Genomewise. *Genome*
352 *Research* **14**, 988-995, doi:10.1101/gr.1865504 (2004).
- 353 28 Bairoch, A. & Apweiler, R. The SWISS-PROT protein sequence database and its
354 supplement TrEMBL in 2000. *Nucleic Acids Res* **28**, 45-48,
355 doi:10.1093/nar/28.1.45 (2000).
- 356 29 Hunter, S. *et al.* InterPro: the integrative protein signature database. *Nucleic Acids*
357 *Res* **37**, D211-215, doi:10.1093/nar/gkn785 (2009).
- 358 30 Harris, M. A. *et al.* The Gene Ontology (GO) database and informatics resource.
359 *Nucleic Acids Res* **32**, D258-261, doi:10.1093/nar/gkh036 (2004).
- 360 31 Kanehisa, M. & Goto, S. KEGG: kyoto encyclopedia of genes and genomes.
361 *Nucleic Acids Res* **28**, 27-30, doi:10.1093/nar/28.1.27 (2000).
- 362 32 Fan, Z., Yue, B., Zhang, X., Du, L. & Jian, Z. CpGIScan: an ultrafast tool for
363 CpG islands identification from genome sequence. *Current Bioinformatics* **12**,
364 181-184 (2017).
- 365 33 Barazandeh, A., Mohammadabadi, M., Ghaderi-Zefrehei, M. & Nezamabadi-
366 Pour, H. Genome-wide analysis of CpG islands in some livestock genomes and

- 367 their relationship with genomic features. *Czech Journal of Animal Science* **61**,
368 487-495 (2016).
- 369 34 Han, L., Su, B., Li, W. H. & Zhao, Z. CpG island density and its correlations with
370 genomic features in mammalian genomes. *Genome Biology* **9**, R79,
371 doi:10.1186/gb-2008-9-5-r79 (2008).
- 372 35 Wright, S. I., Agrawal, N. & Bureau, T. E. Effects of recombination rate and gene
373 density on transposable element distributions in *Arabidopsis thaliana*. *Genome*
374 *Research* **13**, 1897-1903, doi:10.1101/gr.1281503 (2003).
- 375 36 Li, H. *et al.* TreeFam: a curated database of phylogenetic trees of animal gene
376 families. *Nucleic Acids Research* **34**, D572-D580, doi:10.1093/nar/gkj118 (2006).
- 377 37 Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high
378 throughput. *Nucleic Acids Res* **32**, 1792-1797, doi:10.1093/nar/gkh340 (2004).
- 379 38 Guindon, S. *et al.* New algorithms and methods to estimate maximum-likelihood
380 phylogenies: assessing the performance of PhyML 3.0. *Syst Biol* **59**, 307-321,
381 doi:10.1093/sysbio/syq010 (2010).
- 382 39 Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular*
383 *Biology and Evolution* **24**, 1586-1591, doi:10.1093/molbev/msm088 (2007).
- 384 40 Hedges, S. B., Dudley, J. & Kumar, S. TimeTree: a public knowledge-base of
385 divergence times among organisms. *Bioinformatics* **22**, 2971-2972,
386 doi:10.1093/bioinformatics/btl505 (2006).
- 387 41 Hahn, M. W., Demuth, J. P. & Han, S. G. Accelerated rate of gene gain and loss
388 in primates. *Genetics* **177**, 1941-1949, doi:10.1534/genetics.107.080077 (2007).
- 389 42 Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E.
390 M. BUSCO: assessing genome assembly and annotation completeness with
391 single-copy orthologs. *Bioinformatics* **31**, 3210-3212 (2015).
- 392

393 **Tables and figures**

394 **Table 1 Sequencing results for African pike genome assembly.**

Libraries	Raw data		High-quality data	
	Total bases	Sequencing depth	Total bases	Sequencing depth
	(Gb)	(\times)	(Gb)	(\times)
stLFR	118.6	141.2	60.38	71.88
Nanopore	NA	NA	11.02	13.12
Hi-C	66.86	77.47	65.82	76.25

395

396

397 **Table 2 K-mer analysis for African pike genome.**

K	K-mer Number	K-mer Depth	Genome Size (bp)	Used Bases	Used Reads
17	47,775,728,489	48	995,327,676	57,662,825,000	576,628,250

398

399

400 **Table 3 Assembly statistics for the African pike.**

	stLFR		stLFR+ Nanopore		stLFR+ Nanopore+ HiC	
	Contig	Scaffold	Contig	Scaffold	Contig	Scaffold
Number	52,742	26,897	31,798	26,897	31,612	25,653
Length (bp)	793,428,541	859,230,819	848,265,930	864,075,354	846,590,441	862,056,730
Maximum length (bp)	520,087	20,559,849	2,339,471	20,475,446	2,339,471	34,852,849
Average length (bp)	15,043	31,945	26,676	32,125	26,250	32,805
N50	43,947	5,131,134	352,135	5,146,741	347,381	25,843,955
N90	6,949	9,946	9,400	10,077	9,250	9,853
N rate (%)	0.00	7.65	0.00	1.82	0.00	1.79
GC content (%)	41.30	41.30	41.37	41.37	41.37	41.37

401

402

403

404

405

406 **Table 4 Summary of the assembled 29 chromosomes in the African pike.**

Chromosome	Number of contigs	Length of contigs(bp)
1	174	34,852,849
2	155	33,370,094
3	191	32,320,846
4	163	31,998,173
5	169	31,241,732
6	150	30,610,121
7	231	30,524,172
8	175	29,799,086
9	132	29,196,192
10	157	29,175,684
11	145	27,911,487
12	133	27,279,904
13	104	26,912,730
14	144	26,647,357
15	118	25,843,488
16	137	25,843,955
17	130	25,569,178
18	131	25,245,420
19	122	24,922,465
20	153	24,594,709
21	216	23,062,401
22	116	22,653,434
23	166	22,465,777
24	82	21,478,618
25	168	21,151,346
26	121	19,748,208
27	135	18,283,571
28	48	11,790,282
29	76	8,027,360
Total	4,142	742,520,639 (86.1%)

407

408 **Table 5 Prediction of repetitive sequences in African pike.**

Type	Repeat Size (bp)	% of genome
TRF	42,227,324	4.89
RepeatMasker	140,715,741	16.30
RepeatProteinMask	36,712,306	4.25
<i>De novo</i>	292,061,790	33.81
Total	317,311,789	36.73

409 **Table 6 Repeat annotation of the African pike assembly.**

Type	RepBase TEs		TE Proteins		<i>De novo</i>		Combined TEs	
	Length (bp)	% in genome	Length (bp)	% in genome	Length (bp)	% in genome	Length (bp)	% in genome
DNA	64,714,586	7.49	2,115,720	0.25	106,985,346	12.39	140,841,144	16.30
LINE	49,378,300	5.72	31,164,294	3.61	134,752,407	15.6	153,325,343	17.75
SINE	25,662,188	2.80	0.00	0.00	10,933,947	1.27	34,624,631	4.00
LTR	15,943,766	1.85	3,844,230	0.45	67,207,893	7.78	77,013,770	8.92
Other	17,549	0.002	0.00	0.00	0.00	0.00	17,549	0.002
Unknown	0.00	0.00	0.00	0.00	3,180,101	0.37	3,180,101	0.37
Total	140,715,741	16.30	36,712,306	4.25	268,746,546	31.12	284,106,647	32.89

410

411 **Table 7 Statistics of gene annotations in African pike assembly.**

Gene set		Number	Average transcript length (bp)	Average CDS length (bp)	Average exon per gene	Average exon length (bp)	Average intron length (bp)
<i>De novo</i>	Augustus	23,163	14,343	1,392	7.95	175	1,862
	Genscan	29,084	18,389	1,461	8.06	181	2,398
Homolog	<i>I.punctatus</i>	45,520	27,590	2,021	11.57	175	2,418
	<i>D.rerio</i>	26,485	15,268	25,114	8.68	185	3,060
	<i>P.nattereri</i>	45,066	29,558	1,968	11.22	175	2,699
	<i>P.hyp</i>	37,334	29,021	2,084	11.95	175	2,461
	<i>A.mexicanus</i>	41,543	32,069	1,991	11.54	173	2,855
	<i>C.carpio</i>	55,544	13,199	1,176	6.72	175	2,100
Combined	GLEAN	24,314	15,835	1,712	9.77	175	1,610

412

413 **Table 8 Functional annotations of predicted genes in African pike assembly.**

	Database	Number	Percentage (%)
Total		24,314	100
	InterPro	23,989	96.35
	GO	17,227	70.85
	KEGG	21,824	89.76
	Swissprot	23,012	94.65
	TrEMBL	22,970	94.47
unannotated		315	1.30

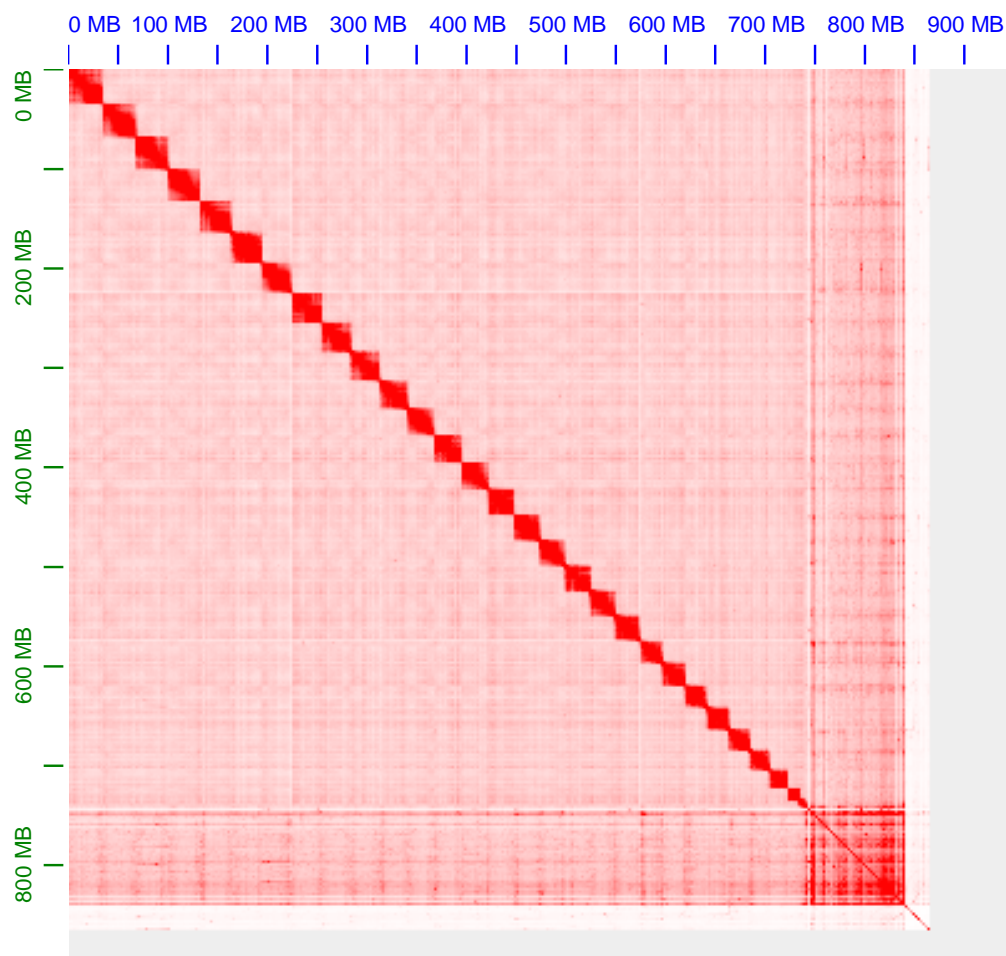
414

415 **Table 9 Statistics of the BUSCO assessment.**

Types of BUSCOs	Gene Set			Assembly		
	Number of actinopterygii (%)	Number of metazoa (%)	Number of vertebrata (%)	Number of actinopterygii (%)	Number of metazoa (%)	Number of vertebrata (%)
Complete BUSCOs	4,241 (92.5%)	943 (96.4%)	2,433 (94.1%)	4,160 (90.7%)	901 (92.1%)	2,531 (90.6%)
Fragmented BUSCOs	222 (4.8%)	26 (2.7%)	110 (4.3%)	252 (5.5%)	14 (1.4%)	169 (6.5%)
Missing BUSCOs	121 (2.7%)	9 (0.9%)	43 (1.6%)	173 (3.8%)	63 (6.5%)	66 (2.6%)
Total BUSCO groups searched	4,584 (100%)	978 (100%)	2,586 (100%)	4,584 (100%)	978 (100%)	2,586 (100%)

416

417 **Figures**



418

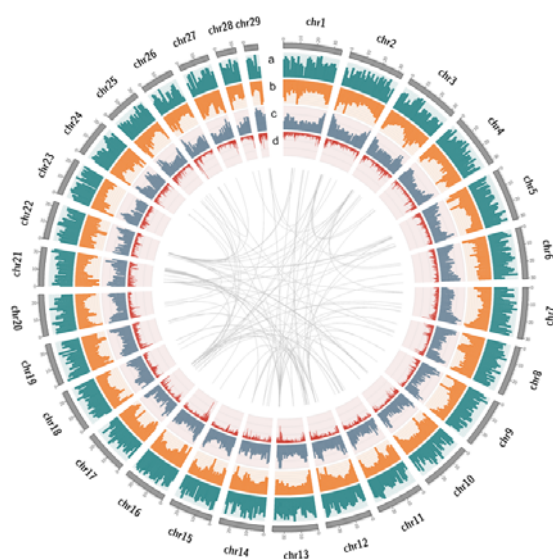
419 **Fig. 1** Hi-C chromosome heat map of African pike genome. Each block represents a Hi-C

420 contact between two genomic loci within a 100kb bin. Darker color represents higher

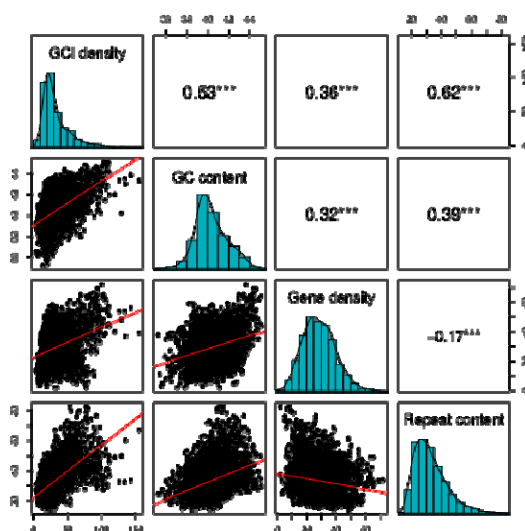
421 contact intensity.

422

423 a



b

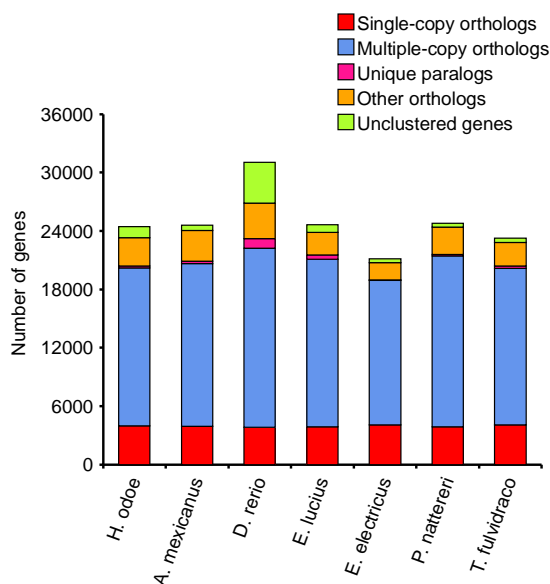


424

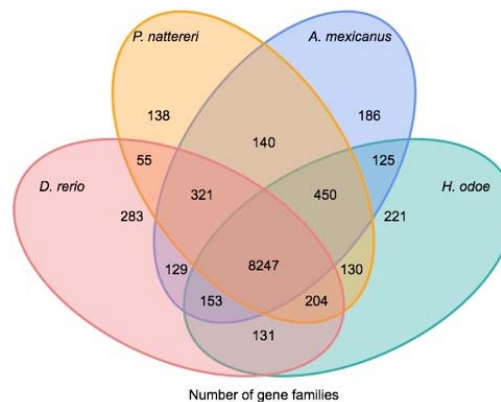
425 **Fig. 2** Genome features of the African pike. (a) The circos plot of 29 chromosomes in
 426 African pike. The tracks from outside to inside are: (1) Gene density, defined as gene
 427 counts per million base pairs. (2) Repeat content, defined as the proportion of repetitive
 428 elements within 1-Mb windows. It was presented in ratio as divided by the highest value.
 429 The axis ranges from 0 to 1. (3) GC content, quantified by the proportion of GC in
 430 unambiguous bases in 1-Mb window. It was presented in ratio as divided by the highest
 431 value. The axis ranges from 0 to 1. (4) CGI content, defined as CGI counts per million
 432 base pairs. The axis ranges from 0 to 200. (b) Correlation matrix among four genome
 433 features. The diagonal presents the distributions by histogram for corresponding genome
 434 features. The lower triangular matrix presents the bivariate scatter plots with a fitted
 435 linear model for each pair of genome features. The upper triangular matrix displays the
 436 Pearson correlation results plus significance level for the corresponding genome features.
 437 Different asterisks represent different significance levels: p -values 0.001 (***), 0.01 (**),
 438 0.05 (*).

439

440 a



b



441

442 **Fig. 3** Comparative genome analysis. (a) Orthologue clustering analysis for African pike

443 and other species. The x-axis displays the seven species and the y-axis presents the gene

444 counts. Red refers to single-copy orthologs, blue refers to multiple-copy orthologs, pink

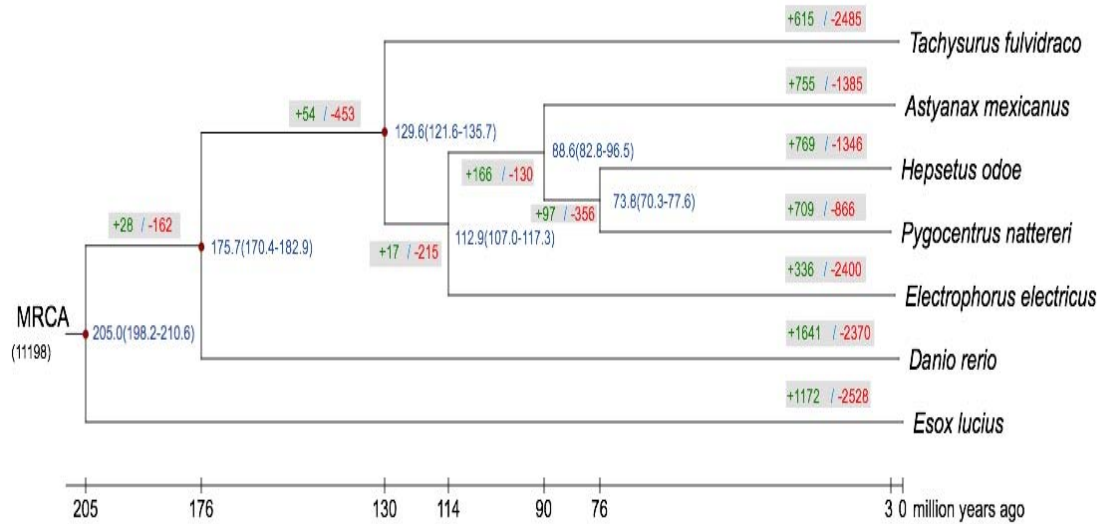
445 refers to unique orthologs for corresponding species, orange stands for other orthologs,

446 and green represents unclustered genes. (b) Venn diagram. Shared and unique gene

447 families among the four species were shown in numbers in corresponding regions.

448

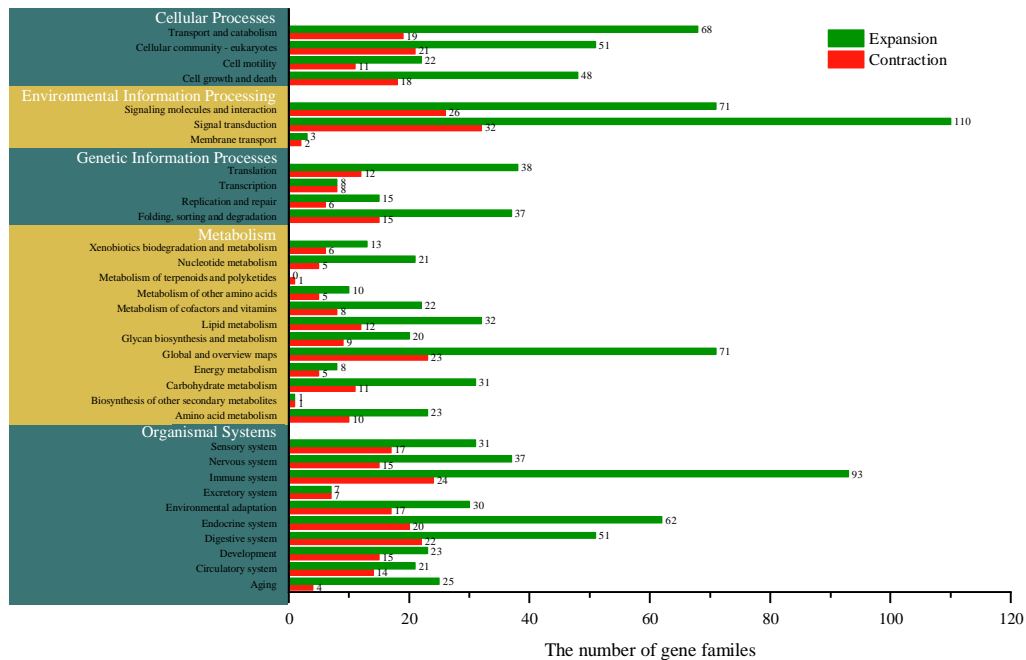
449 a



450

451

b



452

453

454

455

456

457

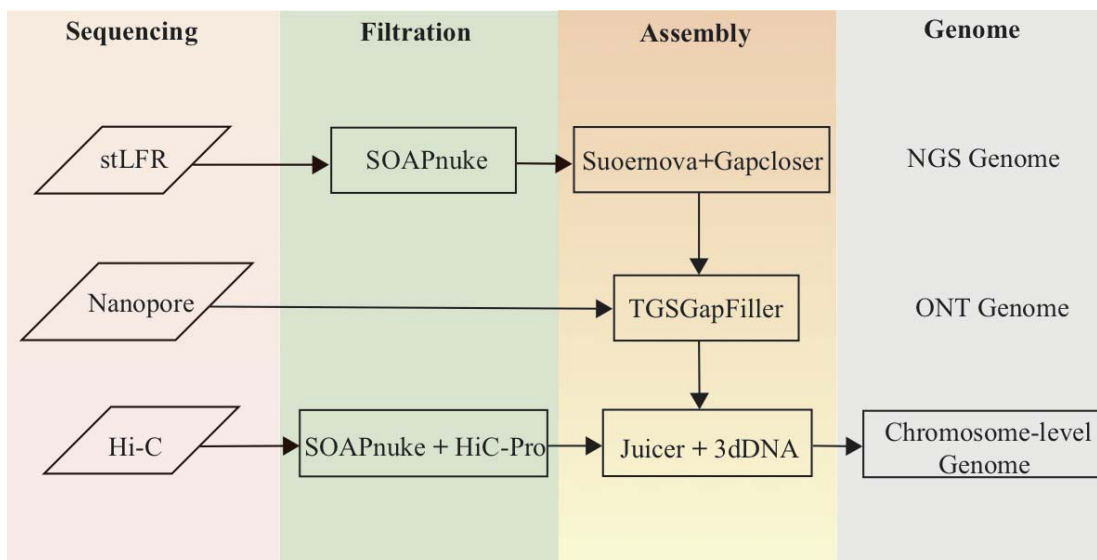
458

459

Fig. 4 Phylogenetic tree, divergence time estimations, and the gene family expansions and contractions for African pike and other species from different fish orders. (a) The phylogenetic tree. Blue values represent the divergence time. Red nodes in the phylogenetic tree represent the reference divergence times. Green and red values represent expanded and contracted gene families for corresponding lineages, respectively. (b) KEGG functional enrichment of the significantly expanded and contracted gene families in the African pike.

460 **Supplementary materials**

461 **Additional file 1: Supplementary figures and tables.**

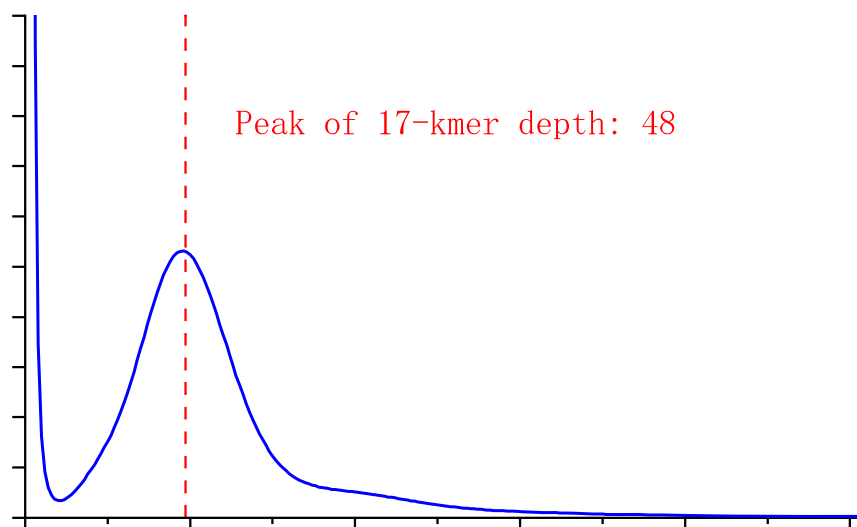


462

463 **Fig. S1** The assembly workflow.

464

465



466

467 **Fig. S2** Distribution of k-mer frequency.

468

469

470 **Table S1 Statistics of gene family clustering.**

Species	Genes number	Genes in families	Unclustered genes	Family number	Unique families	Average genes per family
<i>H.odoe</i>	24,314	23,300	1,014	9,661	103	2.47
<i>A.mexicanus</i>	24,612	24,043	569	9,751	73	2.82
<i>D.rerio</i>	31,056	26,849	4,207	9,523	155	2.82
<i>E. lucius</i>	24,657	23,851	806	9,358	119	2.55
<i>E. electricus</i>	21,162	20,761	401	9,243	18	2.25
<i>P. nattereri</i>	24,796	24,380	416	9,685	57	2.52
<i>T. fulvidraco</i>	23,258	22,796	462	9,389	78	2.43

471

472 **Table S2 Gene family expansion and contraction statistics in African pike genome.**

	Family number	Family number (p <0.05)	Genes Number	KEGG Number
Expansion	769	284	1,548	1,384
Contraction	1,364	81	707	627

473