

Haplotype-resolved telomere-to-telomere assembly of the African catfish (*Clarias gariepinus*) provides insights for semi-terrestrial adaptation of airbreathing catfishes

Julien A. Nguinkal^{1,4*}, Yedomon A. B. Zoclanclounon², Ronald M. Brunner¹, and Tom Goldammer^{1,3*}

¹Research Institute for Farm Animals (FBN), Institute of Genome Biology, Dummerstorf, 18196, Germany

²Department of Crop Science and Biotechnology, Jeonbuk National University, Jeonju, 54896, South Korea

³University of Rostock, Faculty of Agriculture and Environmental Sciences, Rostock, 18059, Germany

⁴Bernhard-Nocht Institute for Tropical Medicine, Department of Infectious Disease Epidemiology, Hamburg, 20359, Germany

*Corresponding authors: Julien A. Nguinkal (julien.nguinkal@bnitm.de), Tom Goldammer (tom.goldammer@uni-rostock.de)

ABSTRACT

Airbreathing catfishes (clariids) are a group of stenohaline freshwater fish that can withstand various environmental conditions and farming practices, including the ability to breathe atmospheric oxygen. This unique ability has allowed them to thrive in semi-terrestrial habitats. However, the underlying genomic and adaptive mechanisms remain poorly investigated. Here, we sequenced the genome of the African catfish *Clarias gariepinus*, one of the most commonly farmed clariids, and generated a gapless telomere-to-telomere (T2T) chromosome-level assembly with high-resolution haplotypes, by integrating long-range sequencing (Hi-C) with PacBio single-molecule (HiFi), Oxford Nanopore, and Illumina sequencing data. The diploid genome assembly yielded 58 contigs with a total length of 969.72 Mb and a contig N50 of 33.71 Mb. We report 25,655 predicted protein-coding genes and 49.94% repetitive elements in the African catfish genome. Our genome assembly provides the first comprehensive gene annotation and haplotype information, such as the male-specific haplotype, enabling us to identify putative genes and molecular mechanisms underlying amphibious traits and terrestrial adaptation of airbreathing catfishes. Several gene families involved in ion transport, osmoregulation, oxidative stress response, and muscle metabolism were expanded or positively selected in clariids, suggesting a potential role in their transition to terrestrial life. The reported findings expand our understanding of the genomic mechanisms underpinning the resilience and adaptive mechanisms of *C. gariepinus* to adverse ecological conditions. They will serve as a valuable resource for future studies in elucidating these unique biological traits in related teleosts and leverage these insights for aquaculture improvement.

Introduction

The *Clariidae* family, commonly referred to as airbreathing catfish, constitutes a group of freshwater fishes that can thrive out of water for extended periods of time by breathing oxygen from the atmosphere^{1,2}. Some of these facultative air breathers have adapted to terrestrial life by developing the ability to walk on land in sinuous movements using their pectoral fins and a protective mucous layer, helping them to retain moisture^{3,4}. These traits enable them to survive in environments with low oxygen levels or stagnant water, such as mangrove swamps, muddy water or flooded forests, which expand their access to new habitats and food sources^{2,5}. To survive in such environments with changing oxygen tensions, this group of fish has developed a bimodal gas exchange capacity in which the gill extracts oxygen from water and the accessory respiratory organ extracts it from the air. Their accessory airbreathing organ (ABO) consists of a paired supra-branchial chamber in the gill cavity. This adaptation of clariids to semi-terrestrial environments (amphibious traits) is however uncommon among bony fish, as only about 11 distantly related fish genera (out of ~2,935)⁶ are considered amphibious with the ability of bimodal respiration³. These independent adaptations and traits diversifications are excellent examples of convergent evolution in teleosts. According to FishBase resources (<https://www.fishbase.se/search.php>), the *Clariidae* family comprises 16 genera and 116 species, with clariids being the most widespread and diverse group with more than 32 recognized species. Many clariids are well-established aquaculture species, including the African catfish (*Clarias gariepinus*, Burchell, 1822), one of Africa's most

31 promising endemic aquaculture fish⁷.

32 *C. gariepinus* is found primarily throughout Africa, where it was first introduced in aquaculture around mid 1970s. This
33 omnivorous fish is quite resilient due to its ability to cope with extreme environmental conditions, tolerate various land-based
34 farming practices and a large diets spectrum^{8–11}. In addition to its rapid growth, extreme robustness^{11,12}, and relatively high
35 fecundity, *C. gariepinus* can withstand high levels of ultraviolet B (UV-B) radiation and dramatic temperature fluctuations
36 in non-aquatic environments^{13,14}. This ecological flexibility could explain its hardiness and wide geographical distribution.
37 Interspecies hybridization with closely related clariids has been shown to improve *C. gariepinus* environmental tolerance,
38 manipulate sex ratios, and eventually increase growth performance, making it a highly efficient aquaculture fish¹⁵. As a
39 result, the African catfish is considered an excellent biological model for studying amphibious traits (i.e., bimodal breathing)
40 and terrestrial transition^{16–18}. However, current genomics research has primarily focused on phylogenetic and domestication
41 studies^{9,19–21}, as well as on sex-chromosome and karyotype evolution utilizing only a limited panel of molecular markers^{22–24}.
42 *Clarias gariepinus* genome is made up of $2n = 2x = 56$ chromosomes (18 m + 20 sm + 18 st/a)²⁵ with a fundamental
43 number (NF) of 94. Its chromosome system has historically been contentious. Previous findings suggested a XX/XY male
44 heterogametic chromosomal system^{26–29}, while others pointed to a ZZ/ZW female heterogametic sex determination system
45 (SDS)^{30,31}. However, recent research using high-throughput sequencing data have shown that both systems coexist in *C.*
46 *gariepinus*^{22,23}. The coexistence of both SDSs is most likely heavily influenced by environmental and social factors, as
47 well as geographical habitat: the ZZ/ZW system is indicated in African wild ecotypes^{25,30}, XX/XY system is observed in
48 some anthropogenically introduced populations in Europe and China^{27,28,32}, and both systems were evidenced within the
49 same population in Thailand^{22,33}. The lack of genomic resources, including reference genomes, haplotypes information, and
50 expression data, has hampered the validation of these SDSs. Yet, few genomic resources of related clariid species, such as the
51 walking catfish (*Clarias batrachus*)³⁴ and the Indian catfish (*Clarias magur*)³⁵, are publicly available. Despite being only at the
52 scaffold levels and highly fragmented with thousands of gaps, these assemblies provide valuable resources for comparative
53 genomic analyses. However, more high-quality genome data are still needed to advance our understanding of the evolution and
54 adaptation of airbreathing catfish to terrestrial habitats. Gold standard genomes, such as telomere-to-telomere (T2T) phased
55 genomes^{36–40}, could facilitate not only studies on sex-chromosome evolution and allele-specific expression, but also provide
56 promising tools for investigating biological mechanisms that shape the robustness and adaptation of airbreathing catfishes.

57 To gain a better understanding of the adaptive strategies of air-breathing fish, we carried out genome sequencing and
58 assembly of the African catfish using HiFi PacBio and Nanopore Technologies, as well as long-range phasing information from
59 Hi-C. We derived a nearly-complete, gapless, and validated phased T2T reference genome assembly. Through a combination of
60 comparative genomic approaches, several genes, biological pathways and processes that are likely associated with the resilience
61 and the emergence of amphibious traits of *Clariidae*, were identified. Our results provide a new genomic basis for functional
62 validation of the molecular mechanisms underlying clariid resilience and their transition out of water, with potential commercial
63 and ecological implications.

64 **Methods**

65 **Sample Collection and DNA extraction**

66 Tissue samples, including muscle, liver, and gonads, were collected from one adult male (approximately one year old) African
67 catfish in the Experimental Aquaculture Facility of the Research Institute for Farm Animal Biology (Dummerdorf, Germany).
68 Prior to tissue collection, the fish was euthanized by immersing it in an overdose of 2-phenoxyethanol (50 mg/L) for 15 minutes,
69 followed by a bleed cut in the head and posterior spinal cord. Tissue samples were immediately frozen in liquid nitrogen, and
70 stored at -80° C. Following the manufacturers' standard protocols, we performed genomic DNA extraction using the DNeasy
71 Blood & Tissue Kit (Qiagen), and libraries preparation strategies specific to the sequencing technologies used in this study.

72 **Libraries preparation and genome sequencing**

73 Genomic DNA (gDNA) sequencing data were generated by different platforms, including Oxford Nanopore (ONT) long reads,
74 PacBio high-fidelity (HiFi) reads, Illumina paired-end reads, and paired-end Hi-C reads (Figure 1a). Illumina short-insert (450
75 bp) libraries were prepared from liver tissues using an Illumina TruSeq Nano DNA Library Prep Kit and paired-end (PE150)
76 sequenced on the Illumina Novaseq 6000 sequencing platform (Illumina, Inc., San Diego, CA, USA). We used gonad tissues
77 for ONT PromethION library preparation and sequencing, following the manufacturer's (Oxford Nanopore Technologies,
78 Oxford, UK) guidelines. In addition, we sequenced a single flow cell on the PromethION instrument, yielding 84 Gb of data
79 and a sequencing depth of around 80 \times , with a maximum read length of 330 kb and an N50 of 32 kb. Liver and muscle tissues
80 were pooled for HiFi library preparation and sequenced on the PacBio Sequel IIe sequencing platform (Pacific Biosciences of
81 California, Inc.). In total, we sequenced four SMRT cells, yielding around eight million CCS reads (141 Gb of data) with an
82 N50 of 16 kb and average base call accuracy greater than 99.7%. A Hi-C library was generated using the Arima-HiC kit and
83 following its standard workflow (Arima Genomics, San Diego, CA, USA). All sampled tissues were pooled and then sequenced

84 paired-end (PE150) on an Illumina HiSeq X platform, yielding 182 million read pairs, corresponding to approximately 24×
85 coverage of the genome. An overview on generated whole-genome sequencing data is provided in the **Supplement Table 1**.

86 **Genome survey analysis**

87 To estimate the preliminary properties of the African catfish genome, we performed a genome-wide *k*-mer analysis using the
88 *k*-mer Analysis Toolkit (KAT) (v2.2.0)⁴¹. Briefly, we generated *k*-mer frequency count (*k* = 21) from high-quality genomic
89 HiFi reads using KAT hist function. With the resulting 21-mer histogram, KAT gcp was used to estimate the genome size,
90 heterozygosity rate, repeat content, and 21-mers derived from errors and sequencing bias. We rendered these genomic properties
91 using a custom R script (Figure 1c-d).

92 **Haplotype-resolved chromosome-scale assemblies**

93 Three strategies were used to generate phased assemblies: hifiasm regular mode, HiFi+Hi-C mode, and haplotype-specific
94 HiFi reads obtained through read partitioning. The output assemblies include a primary assembly (Pim), an alternate assembly
95 (Alt), and two haploid assemblies that include haplotype 1 (Hap1) and haplotype 2 (Hap2). The primary assembly is a more
96 contiguous pseudo-haplotype assembly with long alternating stretches of phased blocks that capture both the homozygous
97 regions and a single copy of the heterozygous alleles. Hap1 and Hap2 are phased assemblies that represent the entire diploid
98 genome, consisting of both parental haplotypes. We used the haplotype-resolved assembler hifiasm (v.0.16.1)⁴² in regular
99 mode (i.e., without Hi-C data) with default parameters to build a contig-level primary and alternate assemblies with clean
100 PacBio HiFi reads. Furthermore, a combination of HiFi and PE Hi-C reads was used in hifiasm to generate a set of two
101 haplotype-resolved, phased contig-level (haplotig) assemblies (i.e., hifiasm Hi-C mode). With `purge_dups` (v1.2.6)⁴³, we then
102 identified and removed contigs corresponding to haplotypic duplications, false duplications, sequence overlaps, and repeats. To
103 construct chromosomes-level phased assemblies, Hi-C PE data were aligned to the purged contigs using a slightly modified
104 Arima Genomics mapping pipeline⁴⁴. SALSA2 (v2.3)⁴⁵ was used to perform chromosomes scaffolding in three iterations.

105 Using the phasing information of the haplotype-specific HiFi reads, we generated a set of two haploid assemblies following
106 the workflow described in Garg et al. (2021)³⁷. In brief, all genomic reads generated in this study (Figure 1) were aligned to the
107 unpolished primary assembly generated in hifiasm regular mode, using minimap2 (v2.2.24)⁴² and BWA-MEM (v0.7.17)⁴⁶ for
108 long and short reads, respectively. We then used HiFi alignments to call heterozygous SNPs using NanoCaller (v3.0.0)⁴⁷. With
109 WhatsHap (v1.4)⁴⁸, we phased heterozygous SNPs utilizing inherent phasing information of HiFi, Hi-C, ONT, and Illumina
110 alignments. For each genotype, we extracted haplotype-specific HiFi long reads, which were then assembled independently
111 with hifiasm regular mode (**Figure 1a**). High-quality chromosome-scale phased assemblies including Hap1 and Hap2 were
112 then built using Ragtag (v2.1.0)⁴⁹. Lastly, the mitogenome was assembled using the mitoHiFi (v2.2)⁵⁰ workflow. To check
113 for the presence of putative contaminations, contigs were searched against all Refseq microbial genomes using Kraken2⁵¹.
114 In addition, a megaBLAST search against non-animals chromosome-level assemblies from RefSeq was performed requiring
115 $e\text{-value} \leq 10^{-5}$, and sequence identity $\geq 98\%$. To fill unresolved gaps between contigs in scaffolds, we applied LR_Gapcloser⁵²
116 with clean HiFi reads. The Hi-C contact maps were visually inspected after polishing and iterative gap-filling to detect potential
117 assembly errors. A few obvious misplacements and orientations of large contigs were identified and manually corrected.

118 **Genome assembly quality assessment**

119 Suitable assembly quality metrics were used to assess the overall completeness and accuracy of the *A. catfish* genome assemblies.
120 Benchmarking Universal Single-copy Orthologs (BUSCO) (v5)⁵³ with the actinopterygii dataset and mapping RNA-Seq data
121 from the same species to genome assemblies were conducted to assess gene completeness. The *k*-mer completeness, phasing
122 accuracy, and heterozygosity of the two haplotype assemblies were evaluated by Merqury (v1.3)⁵⁴. For each assembly, the
123 mapping statistics of the raw NGS reads, including Illumina, Hi-C, ONT, and HiFi, were calculated. In terms of completeness,
124 phasing accuracy, and contiguity, the haplotype-resolved assembly with HiFi + Hi-C outperformed all other approaches. Unless
125 otherwise stated, we used this assembly in the various subsequent analyses in this study.

126 **Identification of the putative sex-specific haplotype**

127 To identify the putative paternal haplotype in our assemblies, the full-length nucleotide sequences of two previously identified
128 and validated male-specific DNA markers (Accession numbers: CgaY1: AF332597; CgaY2: AF332598) were obtained from
129 GenBank^{29,55}. Both sequences, 2.6 kb (CgaY1) and 458 bp (CgaY2) in length, were BLASTed against Hap1, Hap2, and Prim
130 assemblies, requiring stringent mapping criteria (identity >98%; queryCoverage >98%).

131 **Repeats annotation**

132 Assemblies were annotated independently to avoid a skewed comparison. The methods described here were used to annotate
133 genes and repeats in both haplotypes and primary assemblies. RepeatModeler (v2.0.3)⁵⁶ was used to analyze and predict repeat
134 sequences, as well as dependencies such as TRF, RECON, and RepeatScout. Using MITE Tracker⁵⁷, we identified miniature

135 inverted-repeat transposable elements (MITEs). GenomeTools⁵⁸ and LTR_Retrieve (v2.9.0)⁵⁹ were used to analyze full-length
136 LTRs. Furthermore, we retrieved all teleost-specific transposable elements (TEs) from FishTEDB⁶⁰, a curated database of TEs
137 identified in complete fish genomes. We used cd-hit (v4.8.1)⁶¹ to cluster repeat elements with identities greater than 98%.
138 Repeatmasker (v4.1.3)⁶² was used to mask the genome with the final custom non-redundant set of repeats. Utilizing the
139 telomere identification toolkit (tidk)⁶³, we scanned *C. gariepinus* genome for terminal telomeric repeats (5'-TTAGGG-3')_n with
140 a minimum length of 270 bp ($n = 45$) in 25 kb windows of chromosomal termini. To be termed 'terminal telomeric repeats', we
141 required the motif (TTAGGG/CCCTAA)_n to exhibit the highest density per 25 kb in the terminal 25 kb windows compared to
142 internal 25 kb windows. All non-terminal telomeric repeats are referred to as internal or interstitial telomeric sequences (ITS).

143 Genes annotation

144 Protein-coding genes were annotated in *C. gariepinus* genome using *ab initio*, homology-based, and transcriptome-based
145 prediction methods. For homology-based prediction, we obtained high-quality protein sequences from UniProt, which were
146 combined with homologous protein sequences from nine closely related catfish species (**Supplement Table 2**). To map these
147 homologous protein sequences to the African catfish genome, we used TBLASTN with an e-value cutoff of 1e-10. We only
148 kept the highest-scoring alignments with a minimum identity score of more than 80%. The top-scoring proteins were then
149 mapped to the assemblies to predict putative gene models using Exonerate (v2.4.0)⁶⁴. The transcript-based gene prediction
150 was carried out using RNA-Seq data from a conspecific *Clarias gariepinus* individual with available RNA-Seq reads in the
151 Sequence Read Archive (SRA) (BioProject-Accession: PRJNA487132). The quality filtered reads were mapped to our A.
152 catfish assemblies using HISAT2 (v2.2.1)⁶⁵ to detect splice junctions, and StringTie2 (v2.2.0)⁶⁶ was then used to assemble
153 transcripts into gene models.

154 Augustus (v3.4.0)⁶⁷, Genscan⁶⁸, GeneMark-EP⁶⁹, and GlimmerHMM⁷⁰ were used for *ab initio* gene prediction, along
155 with RNA-seq transcript evidence. We used RNA-Seq alignments to train Augustus and GlimmerHMM. In GeneMark-EP and
156 Genscan, we used the default settings. We integrated the genes model prediction from the three methods using the funannotate
157 pipeline (v1.8.13)⁷¹ to build a consensus, non-redundant gene set. Finally, the resulting gene set was filtered to remove genes
158 with no start or stop codon, an in-frame stop codon, or a coding sequence (CDS) shorter than 180 nucleotides (nt). Genes
159 with a high similarity (>90%, e-value < 1e-10) to transposable elements were also removed from the final coding genes set.
160 Several classes of non-coding RNA (ncRNA) genes have also been predicted. tRNAscan-SE⁷² with eukaryote parameters
161 was used to predict transfer RNAs (tRNAs). RNAmmer (v2.1)⁷³ was used to identify eukaryotic ribosomal RNA, and the
162 miRDeep2 pipeline⁷⁴ was used to predict putative microRNAs based on homology to eukaryotic mature miRNA sequences in
163 the miRBASE database⁷⁵.

164 Functional annotation of protein-coding genes

165 The functional annotation of protein-coding genes was achieved by using BLAST to align predicted protein sequences to
166 RefSeq non-redundant proteins (NR) and nucleotides (NT), and UniProtKB/Swiss-Prot databases. Eggnog-mapper (v2.1.9)⁷⁶
167 and Interproscan (v5.56-89.0)⁷⁷ were used to query BLAST top hits (query_coverage > 60%, identity_score > 80%) to obtain
168 Gene Ontology (GO) annotations and gene names via ortholog transfer.

169 Orthologs and phylogenetics analyses

170 The annotated genome *C. gariepinus* allowed us to understand its evolution and estimate divergence time within catfish species.
171 We downloaded protein sequences from NCBI of 14 catfish species from six lineages, including *Clariidae* (five species),
172 *Ictaluridae* (two species), *Siluridae* (one species), *Pangasiidae* (three species), *Bagridae* (two species), and *Sisoridae* (one
173 species). **Supplement Table 2** contains extensive meta-information on these species. Throughout this analysis, two *Cyprinidae*
174 species were used as outgroups: the goldfish (*Carassius auratus*) and the common carp (*Cyprinus carpio*).

175 Gene families from the 14 catfish including outgroup species were identified using the OrthoFinder pipeline with default
176 settings, excepted that the *diamond_more_sensitive* flag was set in alignment parameters. In brief, an all-vs. all BLASTP
177 comparison with an e-value threshold of 1×10^{-10} was performed with all proteins and then genes were clustered into
178 orthogroups using the MCL algorithm. The coding sequences of the single-copy orthogroups were aligned with mafft and
179 concatenated into a super gene for each species. The rooted species tree and gene trees were inferred using single-copy
180 orthologs. The MEGA11⁷⁸ program for Linux was used to estimate the divergence times among the species using rapid
181 relaxed-clock methods⁷⁹ and molecular clock data for calibration constraints obtained from the TimeTree database⁸⁰ between
182 the black bullhead (*Ameiurus melas*) and the goonch (*Bagarius yarrell*).

183 Gene families evolutionary analysis

184 The Computational Analysis of Gene Family Evolution (CAFE) analysis was performed with default parameters to estimate the
185 contraction and expansion of gene families for the 14 catfish species mentioned above. In brief, the time-calibrated ultrametric
186 species tree and orthologous gene families were sent to CAFE (v5)⁸¹, and significant (p - value < 0.05) size variance of gene

187 family expansions and contractions were identified using 1000 random samples, and deviated branches were determined using
188 the Viterbi algorithm implemented in CAFE with a branch-specific p-value less than 0.05. A custom bash script was used
189 to identify significant species-specific gene gain or loss in gene families. Finally, we used the KOBAS-i tool⁸² to perform
190 functional enrichment analyses and to identify pathways and GO terms significantly associated with gene families expansion in
191 the airbreathing catfishes examined in this study.

192 **Positive selection analyses**

193 We used the PosiGene pipeline⁸³ to scan genome-wide positive selection among the aforementioned catfishes, detect selective
194 signatures and understand their role in the adaptive mechanisms of amphibious airbreathing catfishes (*Clariidae*). Positive
195 selection in the *Clariidae* branch was scanned using branch-site tests based on one-to-one single-copy orthologs. The yellow
196 catfish (*Tachysurus fulvidraco*) served as an anchor species, while the black bullhead and goonch served as outgroups. The
197 false discovery rate (FDR) threshold for significantly positively selected genes was set less than 0.05.

198 **Gene duplication events analysis**

199 We examined ten catfish with chromosomal-level genome assembly to identify different types of gene duplication events that
200 could have shaped their evolution. We identified gene pairs derived from whole-genome (WGD), tandem (TD), proximal
201 (PD), transposed (TRD), or dispersed (DSD) duplications using the workflow described by Qiao et al (2019)⁸⁴ and the
202 DupGen_finder pipeline (https://github.com/qiao-xin/DupGen_finder). For each duplicate gene pair, we
203 calculated the synonymous (Ks) and non-synonymous (Ka) nucleotide substitution rates between the two paralogs using the
204 calculate_Ka_Ks_pipeline^{84,85}.

205 **Results**

206 **Whole genome sequencing**

207 Sequencing and assembly of teleosts genomes is particularly difficult due to inherent heterozygosity, retained ohnologs, and high
208 repeat content. In this study, we used a stepwise data integration and assembly validation approach with four complementary
209 NGS technologies to generate a T2T and haplotype-resolved assembly of the African catfish. We sequenced tissues from a
210 male *C. gariepinus* specimen (**Figure 1b**) using Illumina PE reads (~82×), PacBio's HiFi reads (~118×), Oxford nanopore
211 reads (~80×), and Hi-C library sequencing data (~24×) (**Figure 1a, Supplement Table 1**).

212 To conduct genome survey analysis, we used 120 Gb high-quality HiFi data. The *k*-mer analysis (*k* = 21) revealed an
213 estimated genome size of ~980 Mbp, a relatively high heterozygosity rate of 2.12%, and the expected repetitive sequences
214 accounted for approximately 46% of the entire genome (**Figure 1c-d**). The *k*-mer spectra histogram illustrates the high
215 heterozygosity between both haplotypes, with homozygous regions consisting mostly of 2-copy *k*-mers and heterozygous
216 regions consisting mostly of 1-copy *k*-mers, as expected from a diploid genome (**Figure 1c**).

217 The *C. gariepinus* genome was *de novo* assembled using three methods: the standard HiFi-only mode, HiFi+Hi-C mode,
218 and reads partitioning using SNPs phasing information from HiFi, Illumina PE, Hi-C and ONT sequencing data. Except for the
219 HiFi-only mode, which produced a partially phased assembly consisting of a collapsed primary assembly and an incomplete
220 and fragmented alternate assembly, we benchmarked the contiguity and phasing accuracy of haplotigs from both HiFi+Hi-C
221 and reads partitioning approaches. Although there was only a slight difference in assembly contiguity and structural accuracy
222 between the two methods, the assembly obtained with HiFi+Hi-C reached a slightly better accuracy (**Supplement Table 3**).
223 Here, we present the HiFi+Hi-C assembly, which has been extensively validated and used as a reference assembly for the
224 various analyses performed throughout this study.

225 **Haplotype-phased T2T chromosome-scale assembly of the African catfish genome**

226 Following QC filtering and duplicates removal, the initial phased contig-level assembly yielded 58, 142, and 212 sequences,
227 with contigs N50 values of 33.71 Mb, 32.12 Mb, and 19.53 Mb for the Primary, Haplotype-1, and Haplotype-2, respectively.
228 As confirmed later by scaffolding with Hi-C data, more than half (*n* = 34) of the 58 primary contigs represented already entire
229 chromosomes end-to-end, or full-length chromosome arms. After polishing and quality improvement, enhanced fully phased
230 chromosome-scale assemblies were obtained by scaffolding contigs into 28 chromosomes and filling most of the gaps. The
231 chromosomes in the Primary assembly (Prim) were sorted and numbered in order of decreasing physical size. Synteny mapping
232 to Prim was used to number chromosomes of Haplotype-1 (Hap1) and Haplotype-2 (Hap2). The chromosome sizes range from
233 52 Mbp (chr1) to 21 Mbp (chr28), with a median length of 32.3 Mbp. It is possible that the high heterozygosity rate (>2%) of
234 the African catfish genome has facilitated this successful haplotype separation, as it has previously been shown that higher
235 heterozygosity rate aids efficient genome unzipping⁵⁴.

236 Approximately 99% of the assembled genome is spanned by the 28 chromosomes of the Primary assembly, which have
237 no gaps, whereas Hap1 and Hap2 contained only 0.01% and 1.44% unresolved nucleotides (gaps), respectively, mostly in

238 repeat-rich genomic regions. Hi-C analysis identified four chimeric contigs, which were manually examined and corrected. The
239 final haplotype-resolved assembly size for Prim, Hap1 and Hap2 is 969.72 Mb, 972.60 Mb, and 954.24 Mb, respectively. Only
240 Hap2 dramatically increased the N50 metric from 19 Mb to more than 33 Mb at the scaffold-level (**Table 1**). Chromosome-wide
241 analysis of telomeric repeats captured the terminal and tandemly repeated motif ($TTAGGG/CCCTAA$)_n at both chromosomal
242 termini (first and last 25 kbp window) in 21 of 28 *C. gariepinus* chromosomes (**Figure 2a**). Terminal telomeric repeats captured
243 in the first and last 25 kbp windows range in length from 300 bp to 14 kbp, with an average length of 4.5 kbp (**Supplement**
244 **Table 4**). Extending the search window to 1 Mbp did not result in significantly larger copy number of terminal telomeric
245 repeat. Terminal 25-kbp windows had significantly ($p_{\text{adjust}} < 0.01$) larger telomere sizes and densities per kbp than terminal
246 1 Mbp windows (**Figure 2b-c**). This result suggests that the terminal 25 kbp windows captured the majority of full-length
247 telomeric repeats in our *A. catfish* chromosomes assembly, which is consistent with previous findings indicating that the length
248 of telomeric DNA in fish ranges from 2 to 25 kb^{86–88}. We also identified a few internal telomeric sequences with high copy
249 number ($n > 200$). These interstitial or pericentromeric telomeric sequences (ITS) have been evidenced as relics of genome
250 rearrangements in some vertebrates species (**Figure 2a**).

251 Validation of the male-specific marker CgaY1 (AF332598) on Hap1

252 CgaY1 (AF332598)²⁹, a previously identified male-specific marker in *C. gariepinus*, was mapped to only one chromosome
253 in Haplotype-1 and Primary assemblies (identity > 99.14; query Coverage > 96.5; e-value = 0). We found no significant hits
254 on Haplotype-2. CgaY1 is on chromosome 24 at position chr24:20208252-20208717 (Prim) and chr24:20319406-20319871
255 (Hap1). To confirm the absence of this male-specific marker on Hap2, we extracted its flanking sequence (2kb upstream and
256 downstream) and aligned it to chromosome 24 in both the Hap2 and Prim assemblies. We found a single (> 95%) match on
257 Prim but none on Hap2. Although we cannot conclusively determine the Y/W chromosome from this data, we assume that
258 haplotype-1 assembly corresponds most likely to the male-specific haplotype. The Genbank accession numbers for the Primary,
259 Haplotype-1, and Haplotype-2 assemblies are GCA_024256425.1, GCA_024256435.1, and GCA_024256465.1, respectively.

260 Genome structural and functional annotation

261 Integrating ab initio predictions, proteins, and RNA-Seq alignments, we independently annotated the primary assembly and
262 both haploid assemblies. In the collapsed diploid assembly, a total of 25,655 protein-coding gene models were predicted. Hap1
263 and Hap2 yielded slightly lower number of predicted genes, with 23,577 and 24,223, respectively (**Table 1**). Approximately
264 200 genes predicted in Prim were completely missing from Hap1 and Hap2 assemblies. The primary assembly consistently
265 resulted in a better functional annotation, which is to be expected given that the diploid assembly includes both haplotypes and
266 maps a more complete representation of the genome structure. Overall, 87.80% of the 73,455 high-quality proteins across
267 the primary assembly and both haplotypes were assigned a function in at least one of the functional databases searched either
268 through sequence homology or orthologs mapping. (**Table 1**).

269 Repetitive sequences made up half (49.94%) of the *C. gariepinus* genome, which roughly corresponded to the estimated
270 repeat content based on *k*-mers analysis. This relatively high repeat content in the African catfish genome is comparable
271 to that found in other catfishes, including *Clarias magur* (43.72%)³⁵, *Clarias macrocephalus* (38.28%)⁸⁹, *Pangasianodon*
272 *hypophthalmus* (42.10%)⁹⁰ and *Hemibagrus wyckioides* (40.12%)⁹¹. Still, this is significantly higher than in *Clarias batrachus*
273 (30.30%)³⁴, which has a smaller genome size (821.85 Mb). Interspersed repeats are the most abundant class of repetitive
274 elements (46%). Retroelements and DNA transposons accounted for only 12 and 6 percent of the repeatome, respectively
275 (**Supplement Table 5**). The distribution of genes and repeats across the chromosomes followed the typical distribution
276 in vertebrate genomes, with higher gene densities in GC-rich regions and lower gene densities in repeat-rich distal and
277 pericentromeric regions (**Figure 3**).

278 Furthermore, we annotated 6,403 full-length ribosomal RNA, 154 microRNA, and 13,536 transfer RNA throughout the
279 African catfish genome. Remarkably, 96% (6150/6406) of the predicted 5S rRNA genes were all found in a single cluster on a
280 2-Mbp region on both chromosome 4 ($n = 2455$) and chromosome 13 ($n = 3725$). Similarly, 84% (21/25) of the predicted 18S
281 rRNA genes were clustered within the first 500 kbp upstream in the terminal telomeric region of chromosome 27 (**Supplement**
282 **Figure 1**). This result is consistent with earlier findings²⁵ in which 5S rDNA was hybridized on a single site on two *C.*
283 *gariepinus* chromosomes and 18S on just one site on a submetacentric (sm) chromosome (**Supplement Figure 1**). The
284 ribosomal 18S DNA probe did, in fact, hybridize with the sub-telomeric/telomeric region of a medium-sized sm chromosomal
285 pair in *C. gariepinus*, which most likely corresponds to the 500 kb telomeric region on chromosome 27 in this study. The 5S
286 rDNA sequences were identified as a single hotspot in two subtelomeric/acrocentric (st/a) chromosome pairs in *C. gariepinus*,
287 which is most likely the 2 Mbp large 5S rRNA genes cluster our study evidenced on chromosome 4 and chromosome 13 (in the
288 regions 16–18 Mbp) (**Supplement Figure 1**).

289 **Assembly assessment and validation**

290 We performed various assessments to validate the high-quality and completeness of our haplotype-phased African catfish
291 genome assembly, including gene completeness, full-length transcript coverage, read mappability rate, phasing accuracy, and
292 genomic k-mer completeness. The BUSCO completeness (99.10%) was comparable between haplotypes and the primary
293 assemblies. Since we missed only than 0.7% of the expected universal orthologs, the gene space spanned by our genome
294 assembly is nearly complete (**Table 1**). Furthermore, approximately 92% of the *C. gariepinus* transcripts could map on our
295 assemblies (> 90% coverage and >90% identity), indicating their high functional completeness. We also mapped genomic reads
296 to our assemblies to assess structural accuracy and found that more than 96.69% of raw PE reads were concordantly aligned.
297 The alignment rate of ONT, HiFi, and Hi-C reads to the primary assembly was 99.91%, 99.95%, and 100%, respectively. The
298 mapping rates to Hap1 and Hap2 were both greater than 99% (**Supplement Table 6**).

299 Merqury was used to validate the assembly qualities by evaluating phasing completeness and accuracy with haplotype-
300 specific *k*-mers. We expected the sets of haplotype-specific *k*-mers to be completely distinct for a perfectly phased assembly,
301 with no mixture of *k*-mers from both haplotypes. Our data shows that Hap1 and Hap2 are orthogonal with only very few
302 haplotype switches and nearly no contamination (**Figure 4a**). Interestingly, homozygous *k*-mers between both haplotypes
303 were ideally shared in the 2-copy peak. In contrast, a substantial amount of haplotype-specific (heterozygous) *k*-mers was
304 distinct in the 1-copy peak in the spectrum copy number plot (**Figure 4b**). The imbalance of *k*-mers specific to each haplotype
305 representing heterozygous alleles is most likely due to significant differences in sex chromosome sizes.

306 In our haplotype-resolved genome assembly, the phased blocks originating from the wrong haplotype were very small
307 and almost entirely absent when plotting them sorted by size (**Figure 4c**). Moreover, the total phased blocks sizes accounted
308 for 97% and 94% of Hap1 and Hap2 assemblies, respectively. Merqury reported N50 phase block sizes of 3.6 Mbp and 5.5
309 Mbp with only 0.28% switch error rate when a maximum of 100 consecutive switches were allowed within a 20 kbp window
310 (**Figure 4d-e**). The collapsed diploid assembly (Prim) recovered 98.32% of the *k*-mers derived from genomic reads, while the
311 haploid assemblies (Hap1 and Hap2) recovered 83.67% and 82.82%, respectively, demonstrating a high genome completeness
312 (**Figure 4e**). The average base-level accuracy in the Prim assembly was roughly QV42, corresponding to an rate of less than
313 0.01%. Hap1 and Hap2 had a slightly lower accuracy than QV40. It should be noted that haplotypes assemblies were not
314 polished to avoid introducing more switch errors and biased homozygosity (**Figure 4f**).

315 Overall, our assembly quality metrics point to a gapless, fully phased, and telomere-to-telomere (T2T) assembly of the
316 African catfish genome. The majority of these metrics meet or exceed the minimum quality standards⁴⁴ of the VGP consortium.
317 Our reported genome, for example, meets the *c.c.P5.Q42.C98* VGP standard, with *c.c.Pc.Q60.C100*⁴⁴ being the highest
318 standard for finished and gapless T2T vertebrate genomes, such as the recently completed gapless human genome sequence⁹².
319 To the best of our knowledge, this assembly is the first T2T, haplotype-resolved, and the most complete *Siluriformes* (catfish)
320 genome assembly published to date (**Supplement Table 2**).

321 **Phylogeny, divergence time and evolutionary history of catfishes**

322 The comparative phylogenomic analyses performed with OrthoFinder assigned 336,681 (94%) of 390,198 genes to 27,587
323 orthogroups shared among catfishes and two outgroup species (common carp and goldfish). A total of 16,281 genes in *C.*
324 *gariepinus* were found to be orthologous between the 14 catfish species, with 378 of them being single-copy orthologs. The
325 alignments of single-copy orthologs were used to infer the species tree and evolutionary divergence time (**Figure 5**). There
326 were 80 orthogroups comprising 840 genes in total unique to all airbreathing catfish species, with 208 genes specific to *C.*
327 *gariepinus* and spanning 80 orthogroups (**Supplement Table 7**). The vast majority of *C. gariepinus*-specific genes were not
328 characterized in functional databases. Though, ten genes belong to the actin family, eight to the peptidase C13 family, and
329 five to the zinc-finger protein family. According to our estimated phylogenetic tree using protein sequences of all homologous
330 single-copy genes, airbreathing catfishes (*Clariidae* clade) split as a monophyletic group around 98 Mya, which is roughly
331 comparable to the divergence time between rodents and humans (96 Mya) (**Figure 5**).

332 The African catfish diverged from clariids last common ancestor (LCA) about 39.3 million years ago, which is consistent
333 with the current understanding of the historical and geographical distribution of the Clariids, with *Clarias gariepinus* being
334 the only clariid species (in our study) native to Africa⁹³. In contrast, the other *Clarias* species are all endemic to Asia. This
335 result is line in with the parphyly hypothesis that was previously put up for the genus *Clarias*⁹⁴. Due to biogeography and
336 adaptive responses to environmental stressors, the African catfish gradually acquired unique traits and features following the
337 split between the African and Asian *Clarias*⁹⁵. Our phylogeny analysis suggests that, the Asian *Clarias* clade underwent its
338 first speciation event about 25 Mya, which is consistent with the age of the fossil records available for these species⁹⁶.

339 **Comparative gene family evolution of airbreathing catfishes**

340 The expansion and contraction of gene families can play an important role in the adaptation of catfish and other organisms
341 to specific environments by enabling the development and expression of beneficial traits while decreasing the expression
342 of less essential ones. Gene expansion and contraction can lead to potentially gain or loss of to phenotypes. To investigate

343 lineage-specific adaptation of *Clarias*, we used CAFE (Computational Analysis of Gene Family Evolution) to estimate gene
344 family expansions and contractions among 27,587 ortholog groups shared by catfishes, including five airbreathing and nine
345 non-airbreathing catfishes (**Methods**).

346 We found 1,429 and 2,547 gene families that are significantly expanded or contracted in airbreathing catfish. Gene families
347 in *Clarias magur* had the most gene expansion events (1,330), while gene families in *Clarias fuscus* exhibited the most
348 contraction (1,209) events (**Figure 5**). We identified 629 novel expanded and 848 contracted gene families in the *Clarias*
349 *gariepinus* genome. The egalitarian nine homolog gene family (EGLN), the rhodopsin (RHO) gene family, the ferretin (FTH)
350 gene family, and the Carboxypeptidase A (CPA) gene family are some examples of expanded gene families in *C. gariepinus* and
351 other *Clarias* sp. (**Figure 6**). These gene families were all thought to be involved in the environmental adaptation of *Clarias*
352 *magur*, a closely related species to *C. gariepinus*³⁵. The EGLN gene family encodes for prolyl hydroxylase enzymes, which are
353 involved in the regulation of hypoxia-inducible factor (HIF). HIF is a protein that plays a key role in the body's response to low
354 oxygen levels, and prolyl hydroxylase enzymes regulate HIF expression. The duplication of the RHO gene has been proposed
355 as a mechanism for the adaptation of tetrapods⁹⁷ and amphibious fishes^{98–100} to terrestrial environments. The expansion of this
356 gene family in *Clarias* may suggest a critical role in their visual system and light adaptation out of water. Finally, FTH proteins
357 have been associated with iron metabolism and are involved in environment-fish-cross-talk^{101,102}.

358 Expanded gene families in *Clarias* are primarily enriched with ion metal binding, apelin signaling, adrenergic signaling
359 in cardiomyocytes, and neuroactive ligand-receptor interaction pathways, to name only a few. Nucleotide binding
360 (GO:0000166), anatomical structure development (GO:0048856), response to stimulus (GO:0050896), and cytoskeletal motor
361 activity (GO:0003774) are some of the significantly overrepresented GO terms associated with expanded gene families in
362 these facultative airbreathing freshwater fishes (**Supplement Figures 2-4**). Overall, gene family expansion in airbreathing
363 catfishes is primarily characterized by the expansion of gene families encoding for ion transporters and enzymes involved in
364 osmoregulation, metabolism, and energy production. The expansion of these gene families may help airbreathing catfishes cope
365 with the challenges of terrestrial life, such as fluctuating oxygen levels and adapting to new energy sources. The expansion of
366 many gene families involved in cytoskeletal motor activity and anatomical structure development may cause adaptive changes
367 in genes expression to promote the development or modification of specialized anatomical structures, such as gills, labyrinth,
368 blood vessels, and muscles, as well as traits required for low-oxygen environments and efficient terrestrial locomotion and
369 survival.

370 Positive selection in airbreathing *Clarias*

371 The genome-wide screening for positive selection in airbreathing catfish detected nine protein-coding genes under selective pressure
372 ($FDR < 0.05$) when compared to non-airbreathing catfishes (**Supplement Table 8**). For example, the 3-hydroxybutyrate
373 dehydrogenase (BDH1), a member of the short-chain dehydrogenases/reductases (SDR) protein family found in airbreathing
374 catfishes, accumulated up to 14 conserved non-synonymous amino acid substitutions (sites) across *Clarias* species but not
375 in non-airbreathing catfishes. SDR enzymes are known to be involved in the metabolism of lipids and regulating energy
376 balance¹⁰³, which could be important for airbreathing catfishes to preserve energy balance when they are moving on land.
377 Additionally, some SDR enzymes are involved in detoxifying harmful compounds such as pollutants and oxidants in terrestrial
378 environments¹⁰⁴, which can help airbreathing catfishes survive in these harsh conditions.

379 Landscape of gene duplications in catfishes

380 Gene duplication is most likely another driver of airbreathing catfish adaptation. This process can result in the evolution of
381 new genes and the expansion of gene families, which contribute to the acquisition of evolutionary novelty. Among the 25,655
382 coding genes in the African catfish genome, there are 13,809 genes derived from diverse gene duplication events. Based on their
383 duplication mode, DupGen_finder (**Methods**) classified duplicated genes into 5 categories: (i) 496 whole-genome duplicates
384 (WGDs, 3.6%), (ii) 1,463 tandem duplicates (TDs, 10.6%), (iii) 572 proximal duplicates (PDs, 4.14%), (iv) 2,970 transposed
385 duplicates (TRDs, 21.5%), and (v) 8,308 dispersed duplicates (DSDs, 60.16%) (**Figure 7a, Supplement Table 9**). We then
386 estimated the rates of synonymous and non-synonymous substitutions (K_s and K_a) for these five gene categories, and tested for
387 selection pressures including positive and purifying selections.

388 The distribution of K_a , K_s , and K_a/K_s among different modes of duplication showed a striking trend, with proximal and
389 tandem duplications having qualitatively higher K_a/K_s ratios than other modes. The K_s values for PD- and TD-derived gene
390 pairs were relatively low (**Figure 7b-d**). This finding implies that recent TDs and PDs in the African catfish have undergone
391 faster sequence divergence than other paralogs. Whole-genome duplications, on the other hand, are more conserved, with
392 much lower K_a/K_s ratios. The three distinct peaks in the K_s distribution graph of WGD-derived duplicates reflect the three
393 rounds of teleost-specific WGD, with no evidence of catfish-specific WGD events. All retained WGDs (100%) and nearly all
394 TRDs (99.93%) paralogs were subjected to purifying selection. Positive selection was significantly detected in PDs (1.34%)
395 and TDs (0.6%) duplicate gene pairs. Gene duplications were also analyzed in non-airbreathing catfish species. A similar
396 trend was observed in K_a/K_s ratios distribution as well as in the relative proportions of duplicated genes under positive

397 or purifying selection in each paralogs' category. In particular, purifying selection was observed in 100% and 99.59% of
398 WGD-derived duplicate genes in the channel catfish (*Ictalurus punctatus*, IPUN) and in the shark catfish (*Pangasianodon*
399 *hypophthalmus*, PHYP), respectively (Figure 7e-f). These insights suggest that most duplicated genes were either eliminated or
400 diverged very fast after the most recent whole genome duplication events in catfishes. The generally demonstrated hypothesis
401 of rediploidization substantiates this assumption: the genome tends to return to a stable diploid state by losing one copy of each
402 duplicated gene through non-functionalization and subfunctionalization^{105,106}.

403 We performed GO enrichment analysis of tandem and proximal duplicates to determine whether the significant selective
404 pressures observed in TDs and PDs drive the evolution of these genes towards biological functions that support the terrestrial
405 adaptation of *Clarias species*. Tandem and proximal-derived duplicates exhibited divergent functional roles although they
406 shared several enriched GO terms involved in immune response, cytoskeletal motor activity, nervous system, and oxygen
407 binding, which are critical for *Clarias* innate immunity, locomotion and adaptation on land (Supplement file 2: Annotation of
408 duplicated genes). In particular, the tandem duplicated mucin-13-like (MUC13) genes are not only under positive selection, but
409 the MUC gene family has also significantly expanded in all five *Clarias* species included in our analysis (Figure 7g), suggesting
410 a beneficial or adaptive role for these catfish species.

411 In summary, these results show that TDs and PDs are substantially involved in the evolutionary mechanisms for adaptation
412 and diversification of airbreathing catfish, as opposed to WGDs and TRDs, which are subjected to strong purifying selection,
413 preventing them from neofunctionalization and subfunctionalization.

414 Discussion

415 Here, we report the first high-quality chromosome-level, haplotype-resolved T2T assembly of the African catfish genome,
416 an economically and ecologically important airbreathing catfish. Leveraging long reads and Hi-C data, we were able to
417 reconstruct the sequences of both haplotypes with total sizes of 954.24 and 972.60 Mbp. Our fully-phased genome assembly
418 exhibited superior quality metrics based on several indicators such as BUSCO, Merqury, phasing accuracy and functional
419 completeness (Figure 4, Table 1). Haplotype-resolved assemblies provide numerous benefits for genomic-based studies of
420 evolution, conservation, and commercial and disease traits. The reported haplotype-resolved genome sequence and annotation
421 provide a powerful tool for enhanced aquaculture and breeding of *C. gariepinus*. It will, for example, aid in sex determination
422 and allow for a better understanding of structural variations, tissue- and haplotype-specific expression. Furthermore, these
423 genomic resources enable more specific investigations of genomic features such as segmental duplications, hybridization, and
424 structural variant hotspots in this and other closely related catfishes^{36,40,107,108}.

425 Most *C. gariepinus* chromosomes assembly are gapless and resolved from T2T (Figure 2). Telomeres are the protective
426 structures that are found at the ends of chromosomes. In teleosts, they consist of a tandemly repeated DNA hexamer ($(TTAGGG)_n$)
427 and proteins that help to protect the ends of the chromosomes from damage and from being recognized as broken DNA. Our
428 study did not only detect both terminal telomeres in 21 of 28 chromosomes, but also several ITS, mainly located at the
429 pericentromeric regions and along the nucleolar organizer regions (NORs). The absence of high-density terminal telomeric
430 signals at both ends of some chromosomes ($n = 7$) is not necessary due to poor assembly of these regions. The telomeres might
431 be lost or shortened gradually on these chromosomes. In fact, the *C. gariepinus* genome consists of nine subtelomeric/acrocentric
432 (st/a) chromosomes. It has been established that st/a chromosomes have a very short p-arm and that the length of their telomeres
433 is often shorter than that of other chromosome types¹⁰⁹. We observed that a few chromosomes without terminal telomeres
434 at both ends exhibit a high copy number of ITS. This suggests that the terminal scaffold is probably misoriented. Though,
435 these ITS may also indicate relics of ancient chromosomal rearrangements in *C. gariepinus*, including centric and tandem
436 chromosome fusion^{110,111}.

437 Gold standards haplotype-resolved assemblies of commercial fish, such as the one presented here, can aid in the design of
438 optimal haplotypes for intra- or interspecies hybridization by avoiding the combination of known incompatibility of alleles.
439 Furthermore, the availability of the two haplotypes of the African catfish is a turning point for modern genomics-assisted
440 breeding strategies for this species. It will ultimately aid in the development of an A. catfish breeding program. Collectively,
441 our T2T phased assemblies of the A. catfish provide the first and most complete view of its genome to date. It paves the way
442 for a variety of applications, including research into the structure and function of telomeres and their role in chromosomal
443 rearrangements and evolution, the loss or fusion of genetic material, and the diversity of karyotypes and sex-chromosome
444 systems in Claridae.

445 Terrestrial adaptation refers to the process by which aquatic species acquire the ability to live or survive on land for
446 an extended period of time. This process is usually driven by genetic, physiological, and behavioural changes triggered
447 by gene family dynamics, gene duplication events, or positive selection^{112,113}. This evolution can involve many processes
448 and mechanisms, such as changes in body structure, including respiratory and circulatory systems, and sensory and nervous
449 systems^{98,114}. Although they acquire certain benefits from the two worlds, bimodal (aerial and aquatic) airbreathing fish
450 face several challenges when adapting to semi-terrestrial habitats. Hypoxic tension, moisture and humidity loss, prolonged

451 exposure to UV radiation, high-temperature fluctuation, locomotion, and exposure to a different spectrum of pathogens, are
452 few typical challenges that are believed to be the driving forces in the genome remodelling and evolution of aquatic species
453 in these habitats^{115–117}. Gene family dynamics (expansion and contraction) is one of the genome remodelling mechanisms
454 that reflect the evolution of organisms' adaptations to new environments¹¹⁸. Our findings show that significantly expanded
455 gene families in *Clarias* sp. are primarily involved in osmoregulation, anatomical structure development, cytoskeletal motor
456 activity, and stimuli responses. The fluctuating temperature on land will have an impact on osmoregulation and homeostasis
457 via biological processes that regulate ion channels, stress response activation, and osmolyte synthesis^{35,119}. Related gene
458 families such as short-chain dehydrogenases/reductases (SDR), Kv channel interacting proteins (KCNIP), Ferritin (FTH) and
459 hypoxia-inducible factor (EGLN) were significantly associated with these biological processes. We predicted these genes to
460 play a crucial role in the evolution of clariids to semi-terrestrial habitats. For example, FTH plays a role in osmoregulation,
461 particularly in response to changing temperature and salinity¹²⁰. FTH also regulates ion channels and transporters involved
462 in osmoregulation and cells' adaptation to hypoxic stress^{121–123}. In addition, the G-protein-coupled receptor (GPCR) gene
463 family is expanded in the *A. catfish*. The adaptation to terrestrial environments requires fish to maintain proper calcium levels
464 in their bodies as they move between aquatic and terrestrial habitats. Maintaining appropriate calcium levels is crucial for fish
465 on land because calcium is involved in various physiological processes, including muscle contraction, nerve signalling, skeletal
466 development and respiration. These processes may result in structural changes such as a well-developed fish musculature, that
467 facilitate efficient support and movement on land³.

468 Besides gene family expansion, gene duplication is another process that is believed to trigger the acquisition of evolutionary
469 novelty. It has been reported that gene duplication contribute to the emergence of amphibious traits, which enhance the terrestrial
470 transition of aquatic species¹²⁴. We have characterized recent gene duplication events in selected catfishes, including *A. catfish*
471 and other non-airbreathing catfishes. Our findings indicated that TD and PD duplicates exhibited a faster rate of evolution
472 than other modes of duplication such as WGD, DSD, and TRD. Several TD and OD derived duplicates in *C. gariepinus* were
473 found to be specifically under positive selection in clariids, implying their importance in the differential adaptation of these
474 fish species to new habitats and lifestyles. Gene duplication contributes to gene dosage by increasing the number of genes
475 (gene expansion) that are useful in the adaptation continuum in response to new niches and environments^{125,126}. We found
476 evidence of positive selection in BDH1 (3-hydroxybutyrate dehydrogenase), a member of the SDR protein family. With up
477 to 14 accumulated non-synonymous substitutions, this tandemly duplicated gene showed an accelerated rate of evolution in
478 airbreathing catfishes. Previous studies on *Clarias magur*¹²⁷ and in terrestrial mammals¹²⁸ have found that few members of the
479 SDR gene family, including BDH1, were significantly upregulated in response to low oxygen levels, stressing their potential
480 role in adapting and surviving in hypoxic environments. This is consistent with the hypothesis that airbreathing in fish evolved
481 as a response to aquatic hypoxia¹²⁹.

482 Overall, these findings suggest that the transition of airbreathing catfish to terrestrial life may rely on a combination of
483 genetic mechanisms such as gene duplication, gene expansion, and positive selection associated with biological processes
484 that shape environmental adaptation. However, it is important to note that the specific roles of the above mentioned genes
485 and biological process in the adaptation of airbreathing catfish remain hypothetical. Though, these predictions lay a solid
486 basis for future studies and further functional validation to fully understand the specific mechanisms that have facilitated the
487 development of additional capabilities for ecological flexibility of airbreathing catfishes. In order to fully understand the drivers
488 underlying the adaptation and evolution of this group of fish to terrestrial or semi-terrestrial habitats, extensive research would
489 be needed to establish causal relationships. Undoubtedly this haplotype-resolved assembly, along with the characterization
490 of potential genes and genetic changes/mechanisms involved in environmental adaption, establish the fundamentals for such
491 future studies. These may include studying gene expression patterns in these fish in response to different environmental factors
492 and performing functional validation of these genes' function. It could also be insightful to compare the genes and pathways
493 known to be involved in the early evolution and adaption of terrestrial vertebrates to the panel of genes and biological processes
494 hypothesized in this study.

495 Conclusions

496 We have deciphered and annotated the African catfish (*C. gariepinus*) genome, an ecologically and commercially important
497 freshwater airbreathing catfish. This T2T chromosome-level assembly, along with both resolved haplotypes, represent a
498 significant advance in our understanding of *C. gariepinus* genomic makeup. Comparative genomics analysis with related
499 catfishes provided critical insights into the evolutionary mechanisms underlying airbreathing catfish's unique terrestrial
500 adaptation, including genes and pathways associated with hypoxia tolerance, locomotion, skeletal muscle development,
501 respiration, osmoregulation, and antioxidant defense. However, to fully uncover the genomic underpinning of these catfishes'
502 transition from aquatic to terrestrial habitats, further research is needed to validate the specific mechanisms by which these
503 unique genetic changes might have contributed to amphibious traits development in Claridae. Furthermore, this work has
504 demonstrated the utility of HiFi data in achieving fully haplotype-resolved genome assemblies. Overall, this study provides a

505 valuable resource for future studies on the genomic mechanisms underlying catfishes' resilience and adaptive mechanisms to
506 adverse ecological conditions. The insights gained could also be leveraged to improve aquaculture practices and enhance the
507 sustainability of catfish farming.

508 **Data availability**

509 All raw high-throughput sequencing data analysed in this project including Illumina PE, Hi-C, HiFi and ONT sequencing
510 reads are available under NCBI BioProject PRJNA818990. Whole genome assemblies and annotations have been deposited
511 at DDBJ/ENA/GenBank under the accessions GCA_024256425.2 (Primary assembly), GCA_024256435.1 (Haplotype-1)
512 and GCA_024256465.1 (Haplotype-2). The version described in this paper is GCA_024256425.2, GCA_024256435.1,
513 GCA_024256465.1. Further primary data and information on research design are provided at Zenodo ([10.5281/zenodo.7760650](https://doi.org/10.5281/zenodo.7760650)).
514

515 **Code availability**

516 Customs scripts and pipelines used in the data analysis and to create figures are available at <https://github.com/bba10g87/catfish-genome>
517

518 **Acknowledgements**

519 We would like to thank Dr. Alexander Rebl for his constructive comments and inputs towards the interpretation of the results.

520 **Funding**

521 This work was funded by the European Maritime and Fisheries Fund (EMFF). EMFF grant: MV-II.1-LM-014

522 **Author contributions**

523 T.G., R.M.B. and J.A.N. conceptualized the project. T.G. acquired funding. R.M.B. and J.A.N. collected and prepared the
524 tissue samples for sequencing. J.A.N., Y.A.B.Z., T.G. and R.M.B. designed the methodology. J.A.N. and Y.A.B.Z. performed
525 whole bioinformatics analyses and developed the figures. J.A.N. wrote the original draft manuscript. T.G., R.M.B., Y.A.B.Z.
526 and J.A.N. contributed to reviewing and editing the manuscript. All authors read and approved the submitted version.

527 **Competing interests**

528 The authors declare no competing interests.

529 **Figures & Tables**

Table 1. Summary of assembly metric of the *Clarias gariepinus* genome, including the primary (Prim), haplotype-1 (Hap1) and haplotype-2 (Hap2).

Category	Quality Metrics	Primary	Haplotype-1	Haplotype-2
General	Total assembly size (Mb)	969.72	972.60	954.24
	GC content	39.0	38.98	38.93
	Repeat content (%)	49.94	50.07	49.29
Continuity	No. Contigs	58	142	212
	Contig N50 (Mb)	33.71	32.12	19.53
	No. Scaffolds	47	119	98
	Scaffolds N50 (Mb)	33.71	34.0	33.18
	Scaffold L50	12	12	12
	Number of gaps	0	180	115
	% Unplaced sequences (Mbp)	1.01 (12.69)	1.70 (16.5)	2.63 (25.12)
	% Gapless length	100	99.99	98.54
Base accuracy	QV	41.86	38.14	39.39
	k-mer completeness (%)	98.32	83.61	81.93
Structural accuracy	Concondantly mapped PE reads (%)	96.75	96.69	97.81
	BUSCO duplicate (%)	1.42	1.47	1.31
	BUSCO missing (%)	0.7	1.26	1.58
	Reliably phased blocks (%)	—	96.87	94.00
	Proteing coding genes	25,655	23,577	24,223
Functional completeness	BUSCO complete (%)	99.10	97.95	97.76
	NR annotation (%)	87.80	86.17	87.00
	Swissprot/Uniprot annotation (%)	68.23	63.12	64.45
	Transcripts alignment rate (%)	95.52	94.61	94.09

References

1. Bevan, D. J. & Kramer, D. L. The respiratory behaviour of an air-breathing catfish, *clarias macrocephalus* (clariidae). *Can. journal zoology* **65**, 348–353 (1987).
2. Haymer, D. S. & Khedkar, G. D. Biology of selected *clarias* catfish species used in aquaculture. *Isr. J. Aquac.* **74**, 1–15 (2022).
3. Sayer, M. D. Adaptations of amphibious fish for surviving life out of water. *Fish Fish.* **6**, 186–211 (2005).
4. Pace, C. & Gibb, A. C. Sustained periodic terrestrial locomotion in air-breathing fishes. *J. fish biology* **84**, 639–660 (2014).
5. Yatuha, J., Kang’ombe, J. & Chapman, L. Diet and feeding habits of the small catfish, *clarias liocephalus* in wetlands of western Uganda. *Afr. J. Ecol.* **51**, 385–392 (2013).
6. Ravi, V. & Venkatesh, B. The divergent genomes of teleosts. *Annu. review animal biosciences* **6**, 47–68 (2018).
7. Skelton, P. H. & Teugels, G. G. A review of the clariid catfishes (siluroidei, clariidae) occurring in southern Africa (1991).
8. Ducarme, C. & Micha, J.-C. Technique de production intensive du poisson chat africain, *clarias gariepinus*. *Tropicultura* **21**, 189–198 (2003).
9. Hildebrand, M.-C., Rebl, A., Nguinkal, J. A., Palm, H. W. & Baßmann, B. Effects of Fe-DTPA on health and welfare of the African catfish *clarias gariepinus* (Burchell, 1822). *Water* **15**, 299 (2023).
10. Dai, W., Wang, X., Guo, Y., Wang, Q. & Ma, J. Growth performance, hematological and biochemical responses of African catfish (*clarias gariepinus*) reared at different stocking densities. *Afr. J. Agric. Res.* **6**, 6177–6182 (2011).
11. Clols-Fuentes, J., Nguinkal, J. A., Unger, P., Kreikemeyer, B. & Palm, H. W. Bacterial community in African catfish (*clarias gariepinus*) recirculation aquaculture systems under different stocking densities. *Front. Mar. Sci.* **10**, 10.
12. Vitule, J. R., Umbria, S. & Aranha, J. Introduction of the African catfish *clarias gariepinus* (Burchell, 1822) into southern Brazil. *Biol. Invasions* **8**, 677–681 (2006).

- 552 **13.** Sayed, A., Abdel-Tawab, H. S., Hakeem, S. S. A. & Mekkawy, I. A. The protective role of quince leaf extract against the
553 adverse impacts of ultraviolet-a radiation on some tissues of clarias gariepinus (burchell, 1822). *J. Photochem. Photobiol.*
554 *B: Biol.* **119**, 9–14 (2013).
- 555 **14.** Weyl, O., Daga, V., Ellender, B. & Vitule, J. A review of clarias gariepinus invasions in brazil and south africa. *J. fish*
556 *biology* **89**, 386–402 (2016).
- 557 **15.** Rahman, M. A. *et al.* Inter-specific hybridization and its potential for aquaculture of fin fishes. *Asian journal Animal*
558 *veterinary Adv.* **8**, 139–153 (2013).
- 559 **16.** Armelin, V. A. *et al.* The baroreflex in aquatic and amphibious teleosts: Does terrestriality represent a significant driving
560 force for the evolution of a more effective baroreflex in vertebrates? *Comp. Biochem. Physiol. Part A: Mol. & Integr.*
561 *Physiol.* **255**, 110916 (2021).
- 562 **17.** Belão, T., Leite, C., Florindo, L., Kalinin, A. & Rantin, F. Cardiorespiratory responses to hypoxia in the african catfish,
563 clarias gariepinus (burchell 1822), an air-breathing fish. *J. Comp. Physiol. B* **181**, 905–916 (2011).
- 564 **18.** Alimba, C. G. & Bakare, A. A. In vivo micronucleus test in the assessment of cytogenotoxicity of landfill leachates in
565 three animal models from various ecological habitats. *Ecotoxicology* **25**, 310–319 (2016).
- 566 **19.** Tiogué, C. T., Nyadjeu, P., Mouokeu, S. R., Tekou, G. & Tchoupou, H. Evaluation of hybridization in two african
567 catfishes (siluriformes, clariidae): Exotic (clarias gariepinus burchell, 1822) and native (clarias jaensis boulenger, 1909)
568 species under controlled hatchery conditions in cameroon. *Adv. Agric.* **2020**, 1–11 (2020).
- 569 **20.** Kánainé Sipos, D. *et al.* Development and characterization of 49 novel microsatellite markers in the african catfish, clarias
570 gariepinus (burchell, 1822). *Mol. Biol. Reports* **46**, 6599–6608 (2019).
- 571 **21.** Li, Z., Wang, X., Chen, C., Gao, J. & Lv, A. Transcriptome profiles in the spleen of african catfish (clarias gariepinus)
572 challenged with aeromonas veronii. *Fish & Shellfish. Immunol.* **86**, 858–867 (2019).
- 573 **22.** Nguyen, D. H. M. *et al.* An investigation of zz/zw and xx/xy sex determination systems in north african catfish (clarias
574 gariepinus). *Front. Genet.* **11**, 562856 (2021).
- 575 **23.** Nguyen, D. H. M. *et al.* Genome-wide snp analysis of hybrid clariid fish reflects the existence of polygenic sex-
576 determination in the lineage. *Front. Genet.* **13**, 80 (2022).
- 577 **24.** Barasa, J. *et al.* High genetic diversity and population differentiation in clarias gariepinus of yala swamp: evidence from
578 mitochondrial dna sequences. *J. fish biology* **89**, 2557–2570 (2016).
- 579 **25.** Maneechot, N. *et al.* Genomic organization of repetitive dnas highlights chromosomal evolution in the genus clarias
580 (clariidae, siluriformes). *Mol. Cytogenet.* **9**, 1–10 (2016).
- 581 **26.** Liu, S. & Yao, Z. Self-fertilization of hermaphrodites of the teleost clarias lazerea after oral administration of 17- α -
582 methyltestosterone and their offspring. *J. Exp. Zool.* **273**, 527–532 (1995).
- 583 **27.** Liu, S., Yao, Z. & Wang, Y. Sex hormone induction of sex reversal in the teleost clarias lazera and evidence for female
584 homogamety and male heterogamety. *J. Exp. Zool.* **276**, 432–438 (1996).
- 585 **28.** Eding, E., Bouwmans, A. & Komen, J. Evidence for a xx/xy sex determining mechanism in the african catfish clarias
586 gariepinus. In *Presentation at the Sixth International Symposium on Genetics in Aquaculture* (Stirling Scotland, UK,
587 1997).
- 588 **29.** Kovács, B., Egedi, S., Bártfai, R. & Orbán, L. Male-specific dna markers from african catfish (clarias gariepinus).
589 *Genetica* **110**, 267–276 (2000).
- 590 **30.** Ozouf-Costaz, C., Teugels, G. & Legendre, M. Karyological analysis of three strains of the african catfish, clarias
591 gariepinus (clariidae), used in aquaculture. *Aquaculture* **87**, 271–277 (1990).
- 592 **31.** Teugels, G. G. The nomenclature of african clarias species used in aquaculture. *Aquaculture* **38**, 373–374 (1984).
- 593 **32.** Viveiros, A., Eding, E. & Komen, J. Effects of 17 α -methyltestosterone on seminal vesicle development and semen
594 release response in the african catfish, clarias gariepinus. *Reproduction* **122**, 817–827 (2001).
- 595 **33.** Teugels, G., Ozouf-costz, C., Legendre, M. & Parrent, M. A karyological analysis of the artificial hybridization between
596 clarias gariepinus (burchell, 1822) and heterobranchus longifilis valenciennes, 1840 (pisces; clariidae). *J. fish biology* **40**,
597 81–86 (1992).
- 598 **34.** Li, N. *et al.* Genome sequence of walking catfish (clarias batrachus) provides insights into terrestrial adaptation. *BMC*
599 *genomics* **19**, 1–16 (2018).

- 600 **35.** Kushwaha, B. *et al.* The genome of walking catfish *clarias magur* (hamilton, 1822) unveils the genetic basis that may
601 have facilitated the development of environmental and terrestrial adaptation systems in air-breathing catfishes. *DNA Res.*
602 **28**, dsaa031 (2021).
- 603 **36.** Low, W. Y. *et al.* Haplotype-resolved genomes provide insights into structural variation and gene content in angus and
604 brahman cattle. *Nat. Commun.* **11**, 1–14 (2020).
- 605 **37.** Garg, S. *et al.* Chromosome-scale, haplotype-resolved assembly of human genomes. *Nat. biotechnology* **39**, 309–312
606 (2021).
- 607 **38.** Xue, L. *et al.* Telomere-to-telomere assembly of a fish y chromosome reveals the origin of a young sex chromosome pair.
608 *Genome biology* **22**, 1–20 (2021).
- 609 **39.** Deng, Y. *et al.* A telomere-to-telomere gap-free reference genome of watermelon and its mutation library provide
610 important resources for gene discovery and breeding. *Mol. plant* **15**, 1268–1284 (2022).
- 611 **40.** Tian, H.-F., Hu, Q., Lu, H.-Y. & Li, Z. Chromosome-scale, haplotype-resolved genome assembly of non-sex-reversal
612 females of swamp eel using high-fidelity long reads and hi-c data. *Front. Genet.* **13** (2022).
- 613 **41.** Mapleson, D., Garcia Accinelli, G., Kettleborough, G., Wright, J. & Clavijo, B. J. Kat: a k-mer analysis toolkit to quality
614 control ngs datasets and genome assemblies. *Bioinformatics* **33**, 574–576 (2017).
- 615 **42.** Cheng, H., Concepcion, G. T., Feng, X., Zhang, H. & Li, H. Haplotype-resolved de novo assembly using phased assembly
616 graphs with hifiasm. *Nat. methods* **18**, 170–175 (2021).
- 617 **43.** Guan, D. *et al.* Identifying and removing haplotypic duplication in primary genome assemblies. *Bioinformatics* **36**,
618 2896–2898 (2020).
- 619 **44.** Rhie, A. *et al.* Towards complete and error-free genome assemblies of all vertebrate species. *Nature* **592**, 737–746 (2021).
- 620 **45.** Ghurye, J. *et al.* Integrating hi-c links with assembly graphs for chromosome-scale assembly. *PLoS computational*
621 *biology* **15**, e1007273 (2019).
- 622 **46.** Li, H. Aligning sequence reads, clone sequences and assembly contigs with bwa-mem. *arXiv: Genomics* (2013).
- 623 **47.** Ahsan, M. U., Liu, Q., Fang, L. & Wang, K. Nanocaller for accurate detection of snps and indels in difficult-to-map
624 regions from long-read sequencing by haplotype-aware deep neural networks. *Genome biology* **22**, 1–33 (2021).
- 625 **48.** Patterson, M. *et al.* Whatshap: weighted haplotype assembly for future-generation sequencing reads. *J. Comput. Biol.* **22**,
626 498–509 (2015).
- 627 **49.** Alonge, M. *et al.* Automated assembly scaffolding elevates a new tomato system for high-throughput genome editing.
628 *BioRxiv* (2021).
- 629 **50.** Uliano-Silva, M., Nunes, J. G. F., Krasheninnikova, K. & McCarthy, S. A. marcelauliano/mitohifi: mitohifi_v2.0,
630 [10.5281/zenodo.5205678](https://doi.org/10.5281/zenodo.5205678) (2021).
- 631 **51.** Wood, D. E., Lu, J. & Langmead, B. Improved metagenomic analysis with kraken 2. *Genome biology* **20**, 1–13 (2019).
- 632 **52.** Xu, G.-C. *et al.* Lr_gapcloser: a tiling path-based gap closer that uses long reads to complete genome assembly.
633 *Gigascience* **8**, giy157 (2019).
- 634 **53.** Manni, M., Berkeley, M. R., Seppey, M., Simão, F. A. & Zdobnov, E. M. BUSCO Update: Novel and Streamlined
635 Workflows along with Broader and Deeper Phylogenetic Coverage for Scoring of Eukaryotic, Prokaryotic, and Viral
636 Genomes. *Mol. Biol. Evol.* **38**, 4647–4654, [10.1093/molbev/msab199](https://doi.org/10.1093/molbev/msab199) (2021).
- 637 **54.** Rhie, A., Walenz, B. P., Koren, S. & Phillippy, A. M. Merqury: reference-free quality, completeness, and phasing
638 assessment for genome assemblies. *Genome biology* **21**, 1–27 (2020).
- 639 **55.** Balogh, R. E. *et al.* Validation of a male-specific dna marker confirms xx/xy-type sex determination in african catfish
640 (*clarias gariepinus*). *bioRxiv* (2022).
- 641 **56.** Flynn, J. M. *et al.* Repeatmodeler2 for automated genomic discovery of transposable element families. *Proc. Natl. Acad.*
642 *Sci.* **117**, 9451–9457 (2020).
- 643 **57.** Crescente, J. M., Zavallo, D., Helguera, M. & Vanzetti, L. S. Mite tracker: an accurate approach to identify miniature
644 inverted-repeat transposable elements in large genomes. *BMC Bioinforma.* **19**, 348 (2018).
- 645 **58.** Gremme, G., Steinbiss, S. & Kurtz, S. Genometools: a comprehensive software library for efficient processing of
646 structured genome annotations. *IEEE/ACM transactions on computational biology bioinformatics* **10**, 645–656 (2013).

- 647 **59.** Ou, S. & Jiang, N. Ltr_retriever: a highly accurate and sensitive program for identification of long terminal repeat
648 retrotransposons. *Plant physiology* **176**, 1410–1422 (2018).
- 649 **60.** Shao, F., Wang, J., Xu, H. & Peng, Z. Fishtedb: a collective database of transposable elements identified in the complete
650 genomes of fish. *Database* **2018** (2018).
- 651 **61.** Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. Cd-hit: accelerated for clustering the next-generation sequencing data.
652 *Bioinformatics* **28**, 3150–3152 (2012).
- 653 **62.** Smit, A. & Green, P. RepeatMasker. <http://www.repeatmasker.org> (2022). Accessed: 2022-05-20.
- 654 **63.** Brown, M. A Telomere Identification toolKit. <https://github.com/tolkit/telomeric-identifier> (2022). Accessed: 2022-08-20.
- 655 **64.** Slater, G. S. & Birney, E. Automated generation of heuristics for biological sequence comparison. *BMC Bioinforma.* **6**,
656 31 (2005).
- 657 **65.** Kim, D., Paggi, J. M., Park, C., Bennett, C. & Salzberg, S. L. Graph-based genome alignment and genotyping with hisat2
658 and hisat-genotype. *Nat. biotechnology* **37**, 907–915 (2019).
- 659 **66.** Shumate, A., Wong, B., Pertea, G. & Pertea, M. Improved transcriptome assembly using a hybrid of long and short reads
660 with stringtie. *PLOS Comput. Biol.* **18**, e1009730 (2022).
- 661 **67.** Stanke, M. *et al.* AUGUSTUS: ab initio prediction of alternative transcripts. *Nucleic Acids Res.* **34**, W435–439 (2006).
- 662 **68.** Burge, C. & Karlin, S. Prediction of complete gene structures in human genomic DNA. *J. Mol. Biol.* **268**, 78–94 (1997).
- 663 **69.** Bruna, T., Lomsadze, A. & Borodovsky, M. Genemark-ep+: eukaryotic gene prediction with self-training in the space of
664 genes and proteins. *NAR genomics bioinformatics* **2**, lqaa026 (2020).
- 665 **70.** Majoros, W. H., Pertea, M. & Salzberg, S. L. Tigrscan and glimmerhmm: two open source ab initio eukaryotic
666 gene-finders. *Bioinformatics* **20**, 2878–2879 (2004).
- 667 **71.** Palmer, J. Funannotate. <https://github.com/nextgenusfs/funannotate> (2022). Accessed: 2022-07-20.
- 668 **72.** Lowe, T. M. & Chan, P. P. tRNAscan-SE On-line: integrating search and context for analysis of transfer RNA genes.
669 *Nucleic Acids Res.* **44**, W54–57 (2016).
- 670 **73.** Lagesen, K. *et al.* RNAmmer: consistent and rapid annotation of ribosomal RNA genes. *Nucleic Acids Res.* **35**, 3100–3108
671 (2007).
- 672 **74.** Friedlander, M. R., Mackowiak, S. D., Li, N., Chen, W. & Rajewsky, N. miRDeep2 accurately identifies known and
673 hundreds of novel microRNA genes in seven animal clades. *Nucleic Acids Res.* **40**, 37–52 (2012).
- 674 **75.** Griffiths-Jones, S., Grocock, R. J., van Dongen, S., Bateman, A. & Enright, A. J. miRBase: microRNA sequences, targets
675 and gene nomenclature. *Nucleic Acids Res.* **34**, D140–144 (2006).
- 676 **76.** Cantalapiedra, C. P., Hernández-Plaza, A., Letunic, I., Bork, P. & Huerta-Cepas, J. eggNOG-mapper v2: functional
677 annotation, orthology assignments, and domain prediction at the metagenomic scale. *Mol. biology evolution* **38**, 5825–
678 5829 (2021).
- 679 **77.** Blum, M. *et al.* The interpro protein families and domains database: 20 years on. *Nucleic acids research* **49**, D344–D354
680 (2021).
- 681 **78.** Tamura, K., Stecher, G. & Kumar, S. Mega11: molecular evolutionary genetics analysis version 11. *Mol. biology*
682 *evolution* **38**, 3022–3027 (2021).
- 683 **79.** Patel, R. & Kumar, S. On estimating evolutionary probabilities of population variants. *BMC Evol. Biol.* **19**, 1–14 (2019).
- 684 **80.** Kumar, S. *et al.* Timetree 5: An expanded resource for species divergence times. *Mol. Biol. Evol.* **39**, msac174 (2022).
- 685 **81.** Mendes, F. K., Vanderpool, D., Fulton, B. & Hahn, M. W. Cafe 5 models variation in evolutionary rates among gene
686 families. *Bioinformatics* **36**, 5516–5518 (2021).
- 687 **82.** Bu, D. *et al.* Kobas-i: intelligent prioritization and exploratory visualization of biological functions for gene enrichment
688 analysis. *Nucleic acids research* **49**, W317–W325 (2021).
- 689 **83.** Sahn, A., Bens, M., Platzer, M. & Szafranski, K. PosiGene: automated and easy-to-use pipeline for genome-wide
690 detection of positively selected genes. *Nucleic Acids Res* **45**, e100 (2017).
- 691 **84.** Qiao, X. *et al.* Gene duplication and evolution in recurring polyploidization–diploidization cycles in plants. *Genome*
692 *biology* **20**, 1–23 (2019).

- 693 **85.** Qiao, X. Funannotate. https://github.com/qiao-xin/Scripts_for_GB/tree/master/calculate_Ka_Ks_pipeline (2022). Ac-
694 cessed: 2022-08-21.
- 695 **86.** Lund, T. C., Glass, T. J., Tolar, J. & Blazar, B. R. Expression of telomerase and telomere length are unaffected by either
696 age or limb regeneration in danio rerio. *PLoS One* **4**, e7688 (2009).
- 697 **87.** Downs, K. P. *et al.* Characterization of telomeres and telomerase expression in xiphophorus. *Comp. Biochem. Physiol.*
698 *Part C: Toxicol. & Pharmacol.* **155**, 89–94 (2012).
- 699 **88.** Ocalewicz, K. Telomeres in fishes. *Cytogenet. genome research* **141**, 114–125 (2013).
- 700 **89.** Hai, D. M. *et al.* A high-quality genome assembly of striped catfish (pangasianodon hypophthalmus) based on highly
701 accurate long-read hifi sequencing data. *Genes* **13**, 923 (2022).
- 702 **90.** Kim, O. T. *et al.* A draft genome of the striped catfish, pangasianodon hypophthalmus, for comparative analysis of genes
703 relevant to development and a resource for aquaculture improvement. *BMC genomics* **19**, 1–16 (2018).
- 704 **91.** Shao, F. *et al.* Chromosome-level genome assembly of the asian red-tail catfish (hemibagrus wyckioides). *Front. genetics*
705 **12** (2021).
- 706 **92.** Nurk, S. *et al.* The complete sequence of a human genome. *Science* **376**, 44–53 (2022).
- 707 **93.** Pouyaud, L. & Paradis, E. The phylogenetic structure of habitat shift and morphological convergence in asian clarias
708 (teleostei, siluriformes: Clariidae). *J. Zool. Syst. Evol. Res.* **47**, 344–356 (2009).
- 709 **94.** Agnese, J.-F. & Teugels, G. Insight into the phylogeny of african clariidae (teleostei, siluriformes): implications for their
710 body shape evolution, biogeography, and taxonomy. *Mol. Phylogenetics Evol.* **36**, 546–553 (2005).
- 711 **95.** Devaere, S., Jansen, G., Adriaens, D. & Weekers, P. Phylogeny of the african representatives of the catfish family clariidae
712 (teleostei, siluriformes) based on a combined analysis: independent evolution towards anguilliformity. *J. Zool. Syst. Evol.*
713 *Res.* **45**, 214–229 (2007).
- 714 **96.** Otero, O. & Gayet, M. Palaeoichthyofaunas from the lower oligocene and miocene of the arabian plate: palaeoecological
715 and palaeobiogeographical implications. *Palaeogeogr. Palaeoclimatol. Palaeoecol.* **165**, 141–169 (2001).
- 716 **97.** Shi, Y. & Yokoyama, S. Molecular analysis of the evolutionary significance of ultraviolet vision in vertebrates. *Proc. Natl.*
717 *Acad. Sci.* **100**, 8308–8313 (2003).
- 718 **98.** You, X. *et al.* Mudskipper genomes provide insights into the terrestrial adaptation of amphibious fishes. *Nat. communica-*
719 *tions* **5**, 1–8 (2014).
- 720 **99.** Collin, S. P. Early evolution of vertebrate photoreception: lessons from lampreys and lungfishes. *Integr. Zool.* **4**, 87–98
721 (2009).
- 722 **100.** Lin, J.-J., Wang, F.-Y., Li, W.-H. & Wang, T.-Y. The rises and falls of opsin genes in 59 ray-finned fish genomes and their
723 implications for environmental adaptation. *Sci. reports* **7**, 1–13 (2017).
- 724 **101.** Xiong, N.-X. *et al.* Ferritin h can counteract inflammatory response in hybrid fish and its parental species after aeromonas
725 hydrophila infection. *Comp. Biochem. Physiol. Part C: Toxicol. & Pharmacol.* **250**, 109174 (2021).
- 726 **102.** Elvitigala, D. A. S. *et al.* Molecular profile and functional characterization of the ferritin h subunit from rock bream
727 (oplegnathus fasciatus), revealing its putative role in host antioxidant and immune defense. *Dev. & Comp. Immunol.* **47**,
728 104–114 (2014).
- 729 **103.** Kavanagh, K., Jörnvall, H., Persson, B. & Oppermann, U. Medium-and short-chain dehydrogenase/reductase gene and
730 protein families. *Cell. Mol. Life Sci.* **65**, 3895–3906 (2008).
- 731 **104.** Contreras, Á. *et al.* A poplar short-chain dehydrogenase reductase plays a potential key role in biphenyl detoxification.
732 *Proc. Natl. Acad. Sci.* **118**, e2103378118 (2021).
- 733 **105.** Parey, E. *et al.* An atlas of fish genome evolution reveals delayed rediploidization following the teleost whole-genome
734 duplication. *Genome Res.* **32**, 1685–1697 (2022).
- 735 **106.** Sémon, M. & Wolfe, K. H. Rearrangement rate following the whole-genome duplication in teleosts. *Mol. biology*
736 *evolution* **24**, 860–867 (2007).
- 737 **107.** Louro, B. *et al.* A haplotype-resolved draft genome of the european sardine (sardina pilchardus). *GigaScience* **8**, giz059
738 (2019).
- 739 **108.** Takeuchi, T. *et al.* A high-quality, haplotype-phased genome reconstruction reveals unexpected haplotype diversity in a
740 pearl oyster. *DNA Res.* **29**, dsac035 (2022).

- 741 **109.** Sánchez-Guillén, R. *et al.* On the origin of robertsonian fusions in nature: evidence of telomere shortening in wild house
742 mice. *J. evolutionary biology* **28**, 241–249 (2015).
- 743 **110.** Vicari, M. R., Bruschi, D. P., Cabral-de Mello, D. C. & Nogaroto, V. Telomere organization and the interstitial telomeric
744 sites involvement in insects and vertebrates chromosome evolution. *Genet. Mol. Biol.* **45** (2022).
- 745 **111.** Meyne, J. *et al.* Distribution of non-telomeric sites of the (ttagg) n telomeric sequence in vertebrate chromosomes.
746 *Chromosoma* **99**, 3–10 (1990).
- 747 **112.** Mable, B. K., Alexandrou, M. A. & Taylor, M. I. Genome duplication in amphibians and fish: an extended synthesis. *J.*
748 *Zool.* **284**, 151–182, <https://doi.org/10.1111/j.1469-7998.2011.00829.x> (2011). <https://zslpublications.onlinelibrary.wiley.com/doi/pdf/10.1111/j.1469-7998.2011.00829.x>.
- 750 **113.** Romero, P. E., Weigand, A. M. & Pfenninger, M. Positive selection on panpulmonate mitogenomes provide new clues on
751 adaptations to terrestrial life. *BMC Evol Biol* **16**, 164 (2016).
- 752 **114.** Lewis, S. V. The morphology of the accessory air-breathing organs of the catfish, *clarias batrachus*: A SEM study. *J. Fish*
753 *Biol.* **14**, 187–191, [10.1111/j.1095-8649.1979.tb03509.x](https://doi.org/10.1111/j.1095-8649.1979.tb03509.x) (1979).
- 754 **115.** Das, A. B. & Ratha, R. K. Physiological adaptive mechanisms of catfish (siluroidei) to environmental changes. *Aquatic*
755 *Living Resour.* **9**, 135–143 (1996).
- 756 **116.** Banerjee, B., Koner, D., Hasan, R., Bhattacharya, S. & Saha, N. Transcriptome analysis reveals novel insights in
757 air-breathing magur catfish (*clarias magur*) in response to high environmental ammonia. *Gene* **703**, 35–49 (2019).
- 758 **117.** Damsgaard, C. *et al.* Evolutionary and cardio-respiratory physiology of air-breathing and amphibious fishes. *Acta Physiol.*
759 **228**, e13406 (2020).
- 760 **118.** Kim, H. *et al.* Gene family expansions in Antarctic winged midge as a strategy for adaptation to cold environments. *Sci*
761 *Rep* **12**, 18263 (2022).
- 762 **119.** Wang, M.-C. & Lin, H.-C. The air-breathing paradise fish (*macropodus opercularis*) differs from aquatic breathers in
763 strategies to maintain energy homeostasis under hypoxic and thermal stresses. *Front. Physiol.* **9**, 1645 (2018).
- 764 **120.** Mudagandur, S. S., Gopalapillay, G., Vijayan, K. K., Shanker, A. & Shanker, C. Effect of salinity stress on gene expression
765 in black tiger shrimp *penaeus monodon*. *Abiotic biotic stress plants-Recent advances future perspectives* 101–120 (2016).
- 766 **121.** Smith, J. J., O'Brien-Ladner, A. R., Kaiser, C. R. & Wesselius, L. J. Effects of hypoxia and nitric oxide on ferritin content
767 of alveolar cells. *J. Lab. Clin. Medicine* **141**, 309–317 (2003).
- 768 **122.** Robach, P. *et al.* Strong iron demand during hypoxia-induced erythropoiesis is associated with down-regulation of
769 iron-related proteins and myoglobin in human skeletal muscle. *Blood, The J. Am. Soc. Hematol.* **109**, 4724–4731 (2007).
- 770 **123.** Siegert, I. *et al.* Ferritin-mediated iron sequestration stabilizes hypoxia-inducible factor-1 α upon I β s activation in the
771 presence of ample oxygen. *Cell reports* **13**, 2048–2055 (2015).
- 772 **124.** Moriyama, Y. & Koshiba-Takeuchi, K. Significance of whole-genome duplications on the emergence of evolutionary
773 novelties. *Briefings functional genomics* **17**, 329–338 (2018).
- 774 **125.** Qian, W. & Zhang, J. Gene dosage and gene duplicability. *Genetics* **179**, 2319–2324 (2008).
- 775 **126.** Glasauer, S. M. & Neuhauss, S. C. Whole-genome duplication in teleost fishes and its evolutionary consequences. *Mol.*
776 *genetics genomics* **289**, 1045–1060 (2014).
- 777 **127.** Mohindra, V. *et al.* Hypoxic stress-responsive genes in air breathing catfish, *clarias magur* (hamilton 1822) and their
778 possible physiological adaptive function. *Fish & shellfish immunology* **59**, 46–56 (2016).
- 779 **128.** Baze, M. M., Schlauch, K. & Hayes, J. P. Gene expression of the liver in response to chronic hypoxia. *Physiol. genomics*
780 **41**, 275–288 (2010).
- 781 **129.** Scott, G. R. *et al.* Air breathing and aquatic gas exchange during hypoxia in armoured catfish. *J. Comp. Physiol. B* **187**,
782 117–133 (2017).

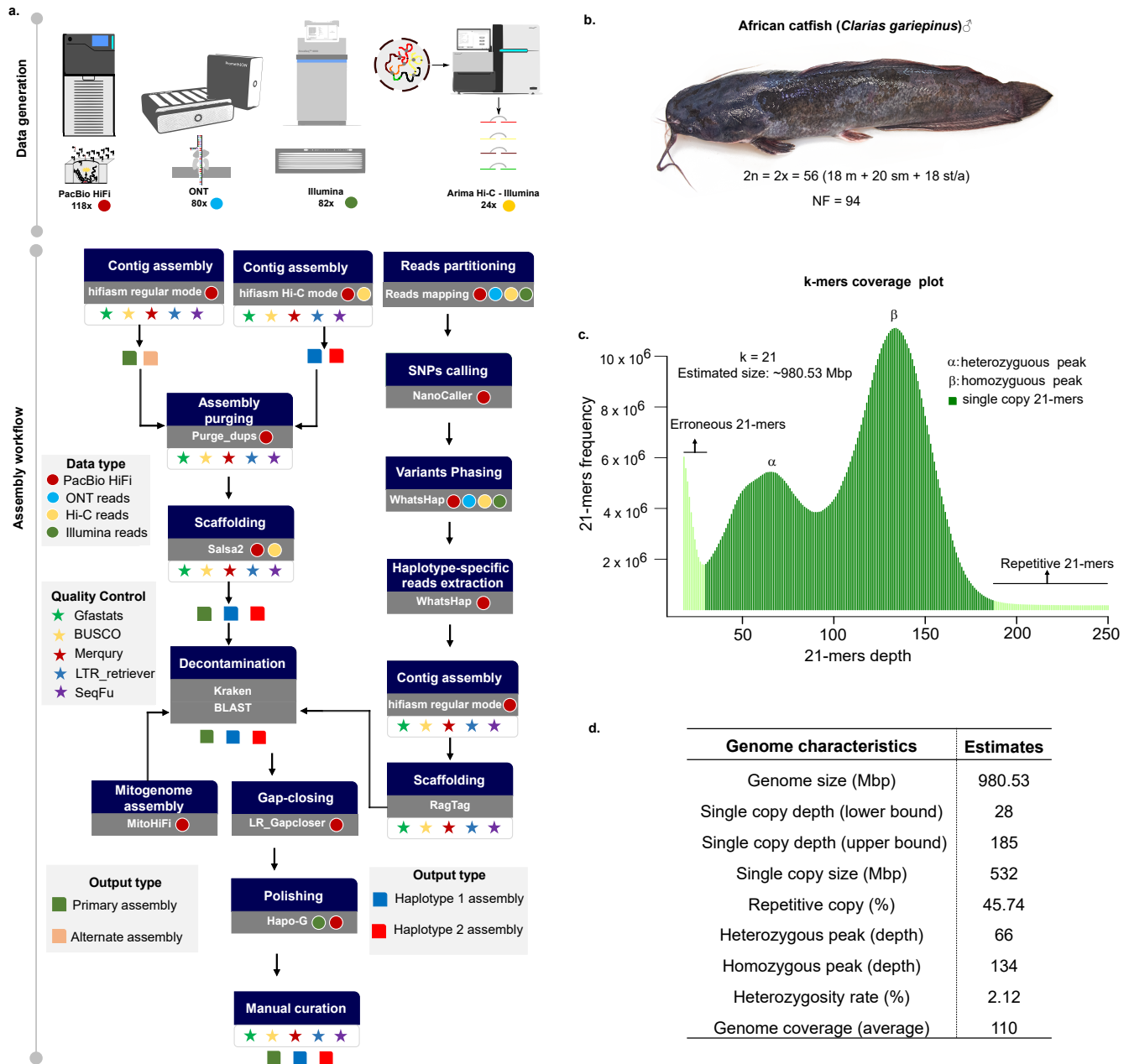
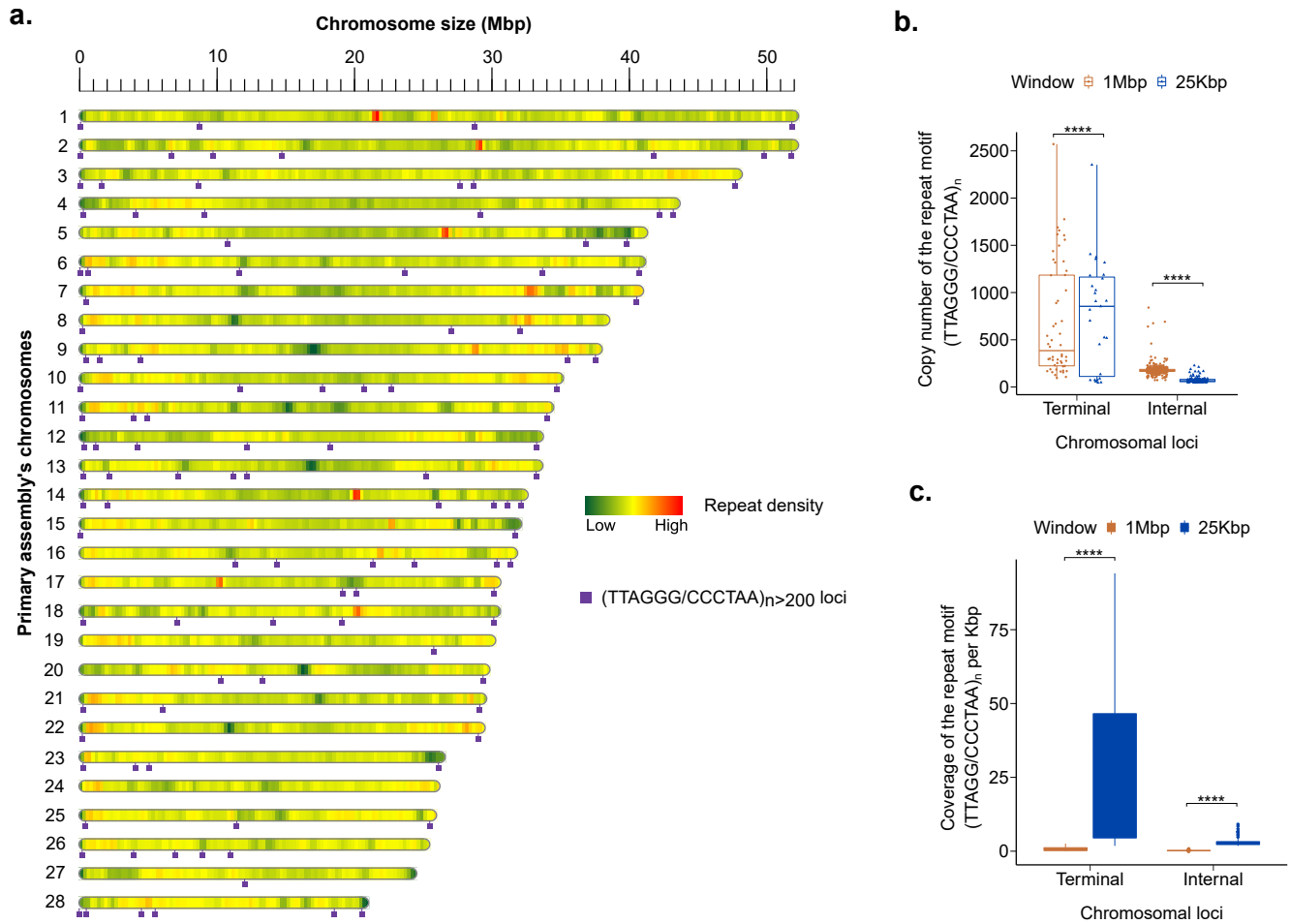


Figure 1. Haplotype-resolved genome assembly workflow of *Clarias gariepinus* and genome survey analysis. **a** The workflow developed to build a haplotype-resolved genome assembly of the the African catfish. Generated genomic sequencing data include Illumina paired-end 150, PacBio’s long high-fidelity (HiFi) reads, Oxford Nanopore (ONT) ultra long reads and Hi-C data. A primary assembly and two haplotype-resolved assemblies were obtained using three assembly modes that combined different data types; **b** The African catfish specimen whose genome was sequenced in this study with the chromosome number for male individuals: A diploid genome with 18 metacentric (m), 20 submetacentric (sm), and 18 subtelomeric/acrocentric (st/a) chromosomes. NF is the fundamental number indicating the total number of chromosome arms; **c** *K*-mer frequency distribution of the diploid genome of the African catfish and its size estimate; **d** Preliminary genome characteristics estimated using *k*-mers analysis.



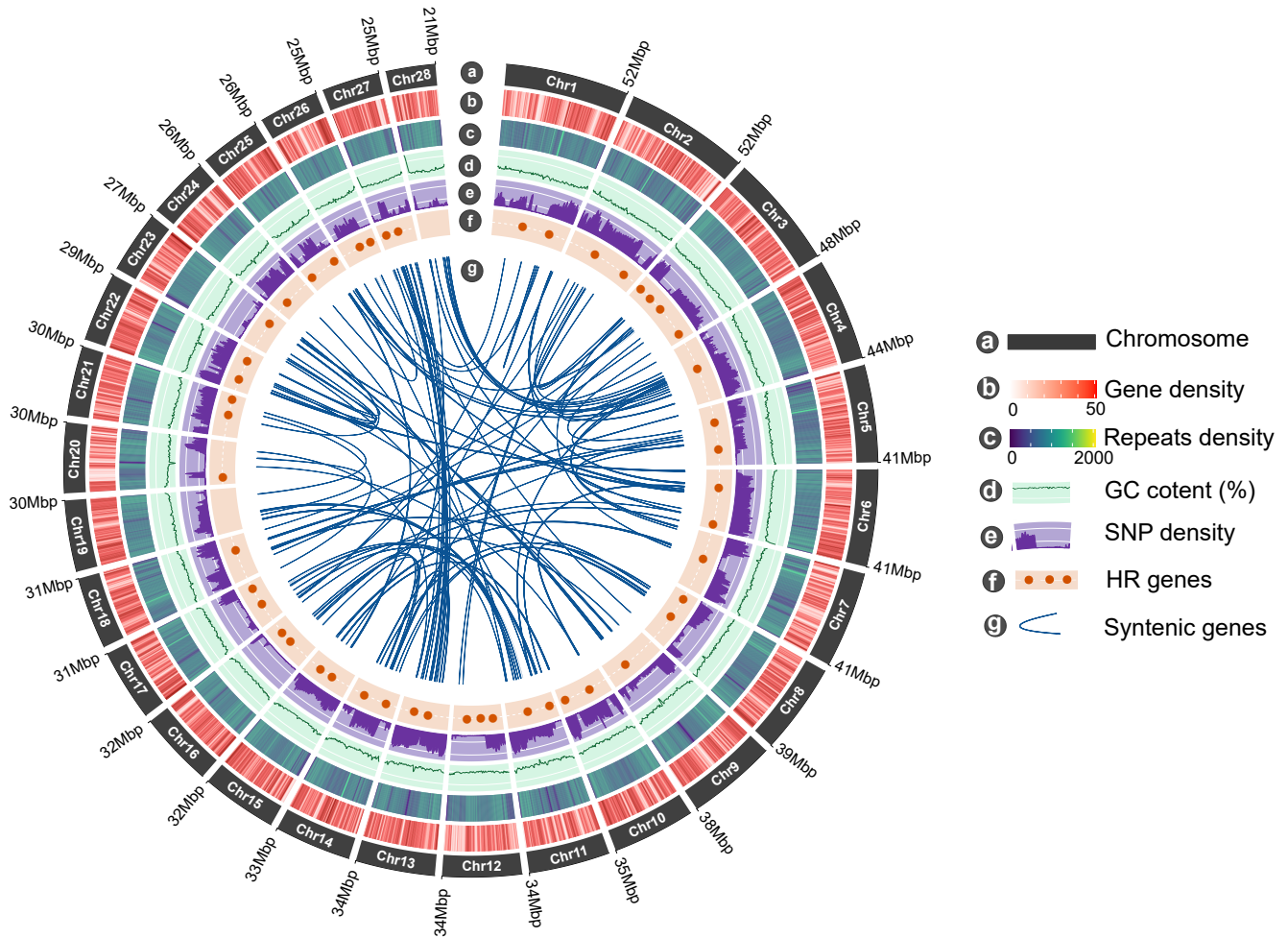


Figure 3. Genomic features of *Clarias gariepinus*. From the outer to the inner circle: **a** Length of the 28 diploid chromosomes (Mb); **b** Chromosome-wide gene density per non-overlapping 500 kb windows; **c** Repeats density in non-overlapping 500 kb windows; **d** GC content; **e** Distribution of heterozygous SNPs density; **f** Chromosomal loci of hypoxia-responsive (HR) genes predicted in the *C. gariepinus* genome; **g** The inner curve lines indicate syntenic gene pairs identified between *C. gariepinus* chromosomes.

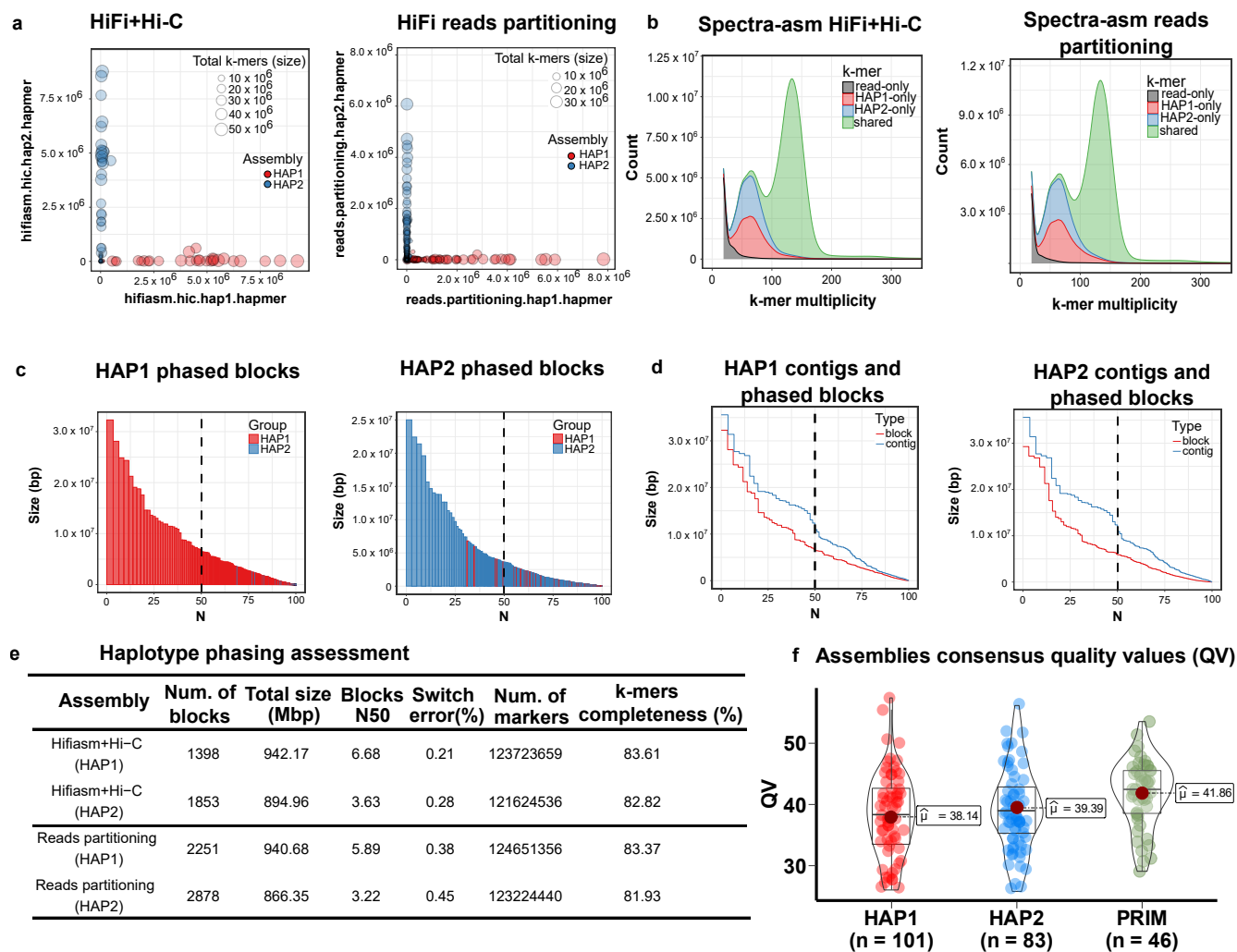


Figure 4. QC plots for evaluating haplotype phasing accuracy, genome contiguity and completeness. **a** Hap-mers blob plot of the Hifi+Hi-C (left) and HiFi reads partitioning assembly (right). Red blobs represent HAP1-specific *k*-mers, while blue blobs are the HAP2-specific *k*-mers. Blob size is proportional to chromosome size. A well-phased assembly should have orthogonal hapmers (e.g. HAP1 and HAP2 lie along axis, respectively). Both assemblies show nearly no haplotypes mixture; **b** Spectra-asm plot of HiFi+Hi-C (left) and Reads partitioning (right) assemblies. The 1-copy *k*-mers representing the heterozygous alleles are specific to each haplotype assembly (HAP1 and HAP2), and the 2-copy *k*-mers, which are only found in the diploid genome, are shared by both assemblies (green). There is no discernible difference between the two assembly approaches. Low-copy *k*-mers (depth < 18) arising from contamination or sequencing errors were removed from the visualization; **c** Phased blocks N* plots of HAP1 (left) and HAP2 (right) assembly, sorted by size. X-axis represents the percentage of the assembly size (*) covered by phased blocks of this size or larger (Y-axis). Blocks from the incorrect haplotype (haplotype switches) are very small and almost entirely absent in the other haplotype. In both haplotypes, more than 75% of the assembly is spanned by phased blocks larger than 1 Mbp; **d** Phase block and contig N* plots showing the relative continuity of HAP1 (left) and HAP2 (right); **e** Statistics for haplotype phasing with switch errors and phased blocks allowing up to 100 switches within 20 kbp; **f** The average consensus quality (QV) distribution for each assembly. Each dot represents a scaffold in the associated assembly.

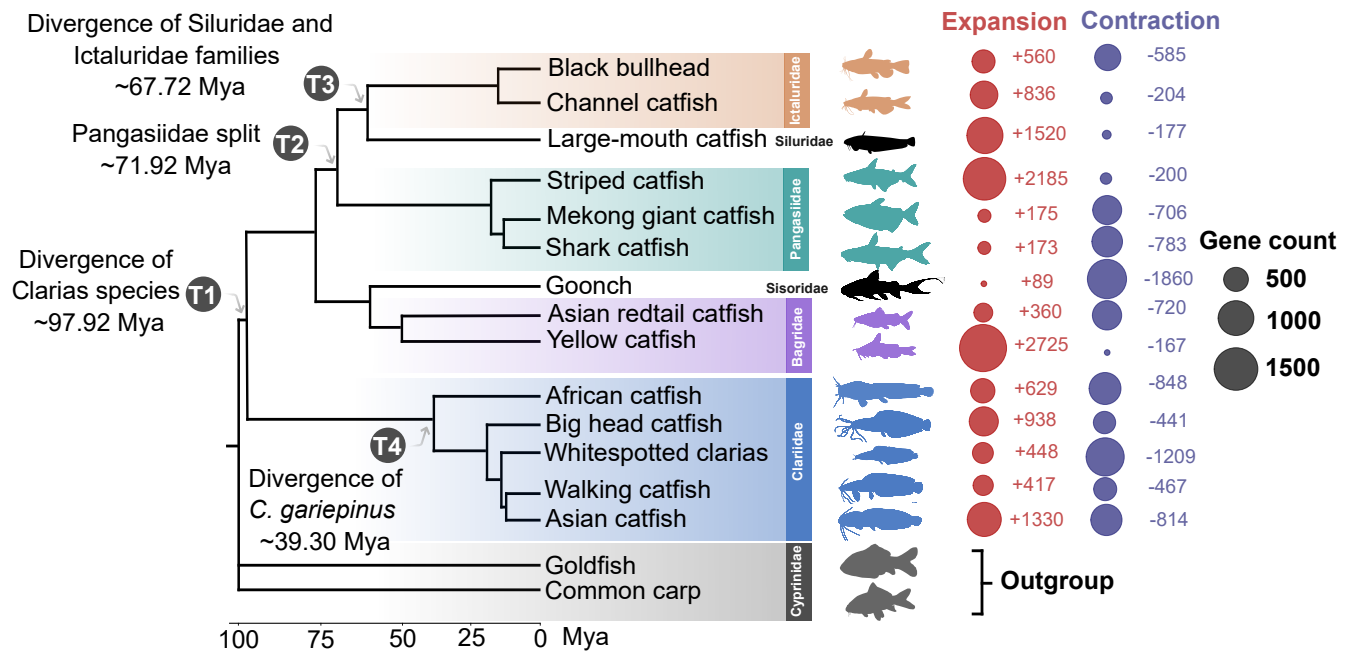


Figure 5. Phylogenomic relationships of major groups of catfishes (Siluriformes). Time-calibrated phylogenetic tree of 14 catfish species based on 1:1 single-copy orthologous proteins. Estimated divergence time as well as the time scale in million years (mya) are shown at the bottom axis. The bubble chart at the right end of the species represents the proportion of gene families that underwent expansion (red) or contraction (blue) in a specific branch. The circle radius is proportional to the number of genes assigned to each category.

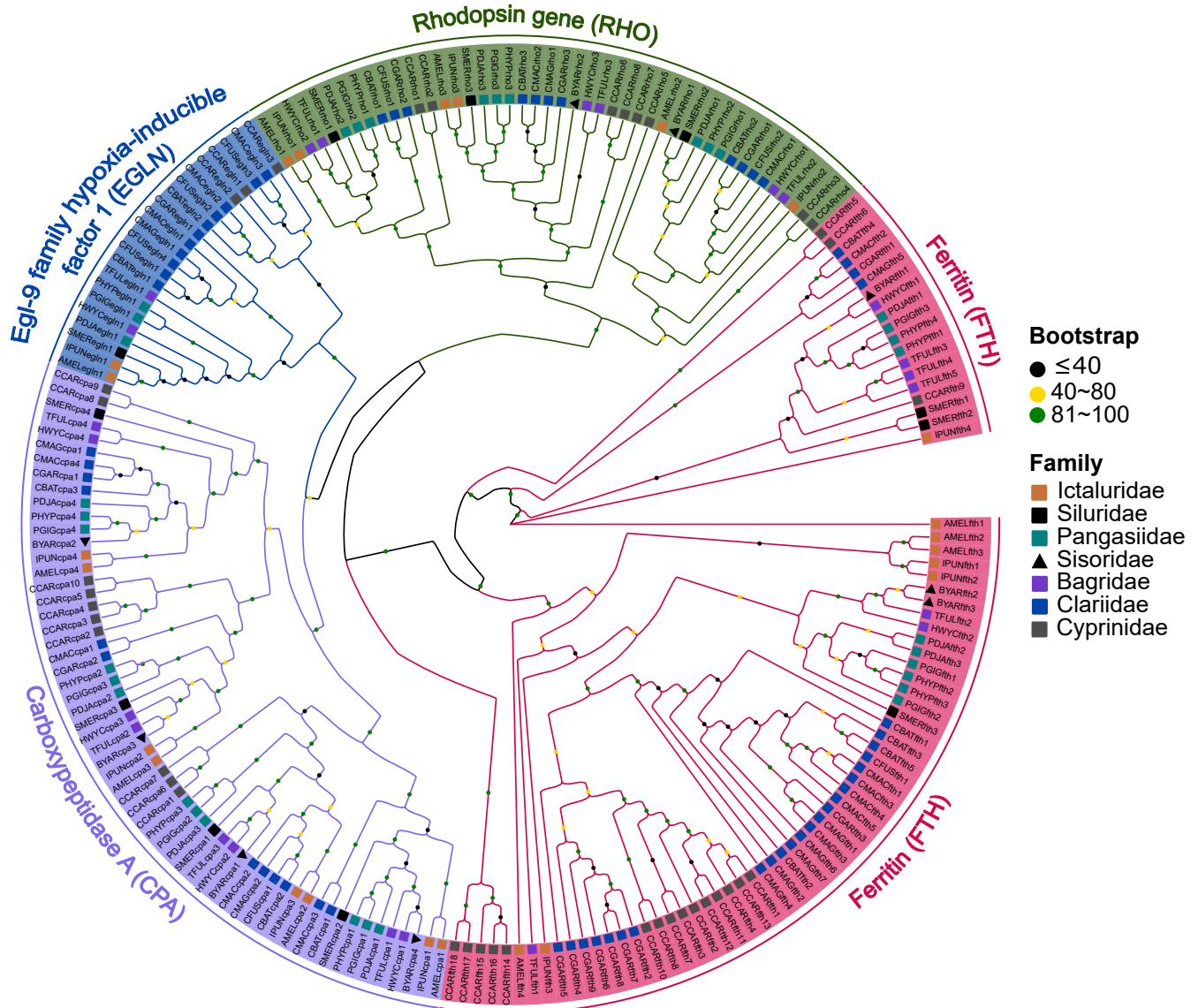


Figure 6. Examples of Clariidae-specific expanded gene families. Maximum likelihood gene tree showing the phylogenetic relationship of four gene families significantly expanded only in *Clariidae* (airbreathing catfishes), but not in non-airbreathing catfishes. Species of the same taxonomic family have the same shape and color. Bootstrap values are indicated with black, yellow and green colors.

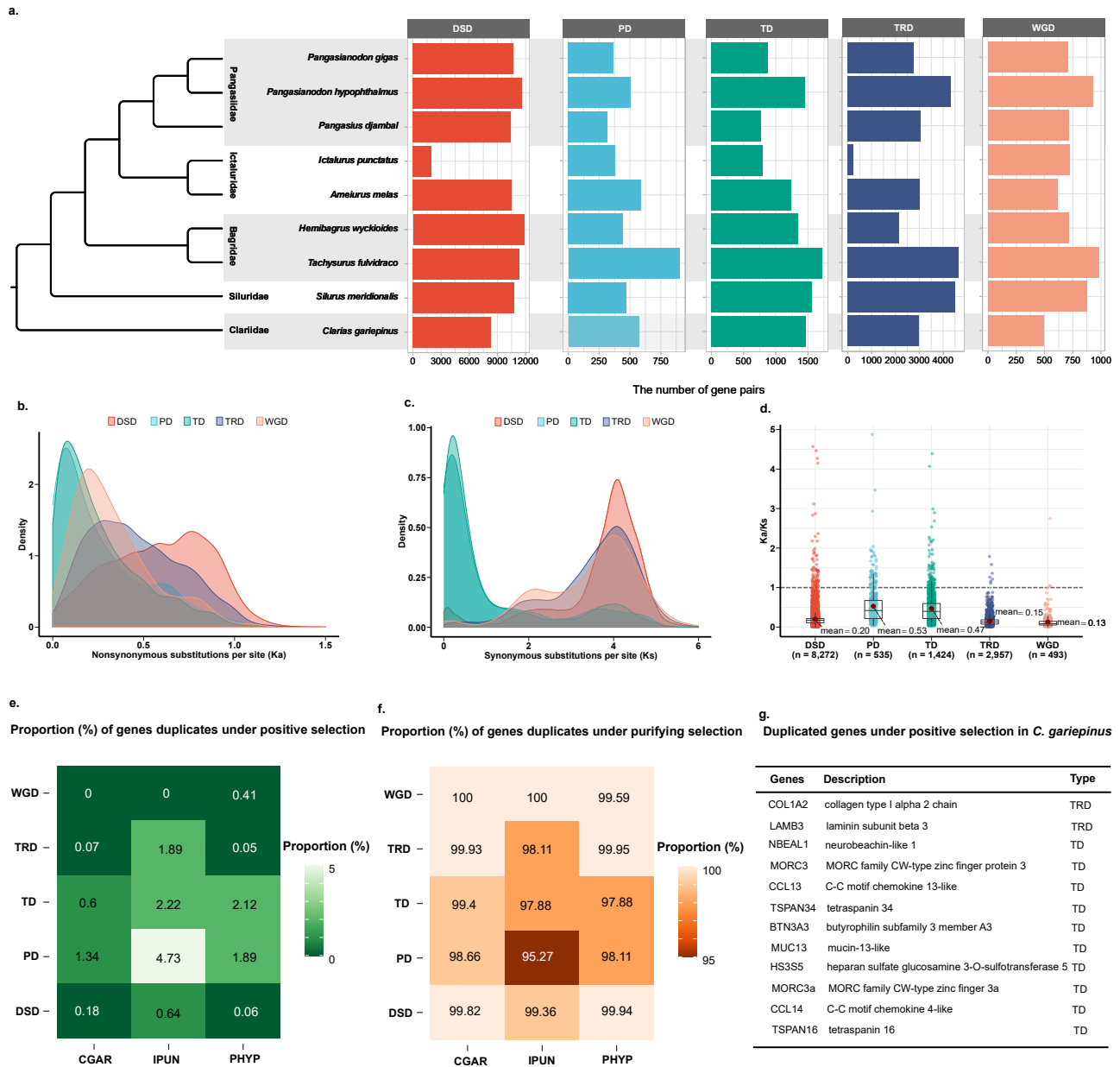


Figure 7. Landscape of gene duplication and positive selection in the A. catfish. **a** The number of gene pairs derived from various duplication modes in representative catfish genomes. DSD dispersed duplication, PD proximal duplication, TD tandem duplication, TRD transposed duplication, and WGD whole-genome duplication are the different types of duplication. It also shows a schematic representation of the phylogeny of the various catfish species used in the study; **b,c** Evolution of gene pairs duplicated by different modes in *A. catfish*. Ka distributions (**b**) and Ks distributions (**c**); **d** The Ka/Ks ratio distributions of gene pairs derived from different modes of duplication in the African catfish; **e** Percentage of genes under positive selection in three catfish lineages; **f** Percentage of genes under purifying selection in three catfish lineages. CGAR: The African catfish (*Clarias gariepinus*), IPUN: The channel catfish, (*Ictalurus punctatus*), PHYP: shark catfish (*Pangasius hypophthalmus*); **g** Duplicated genes in *C. gariepinus* that are positively selected in all clariids examined in this study.