

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11 **Large inversions in Lake Malawi cichlids are associated with**  
12 **habitat preference, lineage, and sex determination**  
13  
14  
15  
16  
17  
18

19 **Nikesh M. Kumar<sup>1</sup>, Taylor L. Cooper<sup>1</sup>, Thomas D. Kocher<sup>2</sup>, J. Todd Streebman<sup>1†</sup>, and Patrick T.**  
20 **McGrath<sup>1,3†</sup>**  
21

22 <sup>1</sup> **School of Biological Sciences, Georgia Institute of Technology, Atlanta, GA, USA**

23 <sup>2</sup> **Department of Biology, University of Maryland, College Park, MD USA**

24 <sup>3</sup> **School of Physics, Georgia Institute of Technology, Atlanta, GA, USA**

25 <sup>†</sup> **Co-correspondence: [patrick.mcgrath@biology.gatech.edu](mailto:patrick.mcgrath@biology.gatech.edu), [todd.streebman@biology.gatech.edu](mailto:todd.streebman@biology.gatech.edu)**  
26

27

28

## 29 **ABSTRACT**

30 *Chromosomal inversions are an important class of genetic variation that link multiple alleles together into*  
31 *a single inherited block that can have important effects on fitness. To study the role of large inversions in*  
32 *the massive evolutionary radiation of Lake Malawi cichlids, we used long-read technologies to identify four*  
33 *single and two tandem inversions that span half of each respective chromosome, and which together*  
34 *encompass over 10% of the genome. Each inversion is fixed in one of the two states within the seven major*  
35 *ecogroups, suggesting they played a role in the separation of the major lake lineages into specific lake*  
36 *habitats. One exception is within the benthic sub-radiation, where both inverted and non-inverted alleles*  
37 *continue to segregate within the group. The evolutionary histories of three of the six inversions suggest they*  
38 *transferred from the pelagic Diplotaxodon group into benthic ancestors at the time the benthic sub-*  
39 *radiation was seeded. The remaining three inversions are found in a subset of benthic species living in deep*  
40 *waters. We show that some of these inversions are used as XY sex-determination systems but are also*  
41 *likely limited to a subset of total lake species. Our work suggests that inversions have been under both*  
42 *sexual and natural selection in Lake Malawi cichlids and that they will be important to understanding how*  
43 *this adaptive radiation evolved.*

## 44 **INTRODUCTION**

45 Large genomic inversions are a particularly interesting type of genetic variation in the evolutionary  
46 process<sup>1-4</sup>. These rearrangements of DNA sequence strongly suppress recombination between inverted  
47 and non-inverted alleles, enabling the capture and accumulation of small genetic variants that are linked  
48 together on each inversion haplotype<sup>5</sup>. These alternative haplotypes largely follow independent  
49 evolutionary paths within a species, accumulating new mutations that can affect phenotypes under  
50 natural and sexual selection. Structural rearrangements can be large, spanning multiple megabases of  
51 DNA containing hundreds of genes<sup>2,6</sup>. They have been shown to play a role in adaptive divergence within a  
52 variety of species<sup>7-12</sup>, in creating alternative reproductive strategies and mating types<sup>13-16</sup>, in evolution of  
53 sex chromosomes<sup>1,17</sup>, and in the formation of pre- and post-zygotic barriers between incipient species<sup>18,19</sup>.

54 When comparing the genomes of extant species, chromosomal inversions are often found as fixed  
55 differences, suggesting they might have an important role in speciation<sup>20-22</sup>. However, the forces  
56 responsible for establishing inversions remain obscure. Do they play a role in adaptation to new ecological  
57 niches exploited by incipient species? Do they carry genetic variants that create prezygotic barriers that  
58 drive speciation? Can they resolve sexual conflict, linking sex-specific alleles to sex determining loci to  
59 create sexually dimorphic species? Perhaps most likely, they play roles in many or all these processes<sup>23,24</sup>.

60 Few vertebrates offer as much potential for investigating evolution and speciation as the *Cichlidae* family  
61 of fish, one of the most species-rich and diverse families of vertebrates, with an estimated 2,000-3,000  
62 species across the globe<sup>25-27</sup>. Evolutionary radiations have occurred at least three times in the African  
63 Great Lakes (Lake Malawi, Lake Tanganyika, and Lake Victoria). In Lake Malawi, more than 800 species of  
64 cichlids evolved in the past 1.2 million years<sup>26,28-30</sup>. Speciation in cichlid fishes involved extremely high  
65 levels of phenotypic divergence, including changes in body shape, jaw structures, and feeding behaviors  
66 related to prey in their ecological niches<sup>25,31</sup>. Additionally, traits that create prezygotic barriers between  
67 species are also extremely diverse. Color patterns are quite dramatic and diverse among species and  
68 include both differences in color and species-specific patterns including bars, stripes, and blotches that  
69 are thought to reinforce barriers between species in sympatry<sup>32</sup>. In some species, mating is seasonal, with  
70 gatherings on 'leks' where males build bowers by scooping and spitting out sand over the course of a week  
71 or more<sup>33</sup>. Females choose mates based upon features of these bowers including shape and size<sup>34</sup>. Most

of the Lake Malawi species can be induced to interbreed in the lab, providing the opportunity to use genetic approaches to study causality of associated changes. Despite the presence of these prezygotic barriers, substantial gene flow has occurred multiple times within the lake<sup>29,35</sup> and hybridization has been proposed to seed radiations within the African Great Lakes<sup>27,36,37</sup>.

Our current understanding of how the Lake Malawi radiation created >800 species is a set of three serial diversifications, starting from a single riverine-like ancestor lineage, that created seven major ecogroups separated primarily by habitat<sup>29</sup>. A pelagic lineage separated first, further diversifying into a deep-water ecogroup (*Diplotaxodon*) and a mid-water ecogroup (*Rhamphochromis*). A muddy/sandy benthic lineage evolved from the ancestral riverine lineage next, which split into three ecogroups (collectively referred to as *benthics*) living either in the water column (*utaka*), over shallow sandy/muddy shores (*shallow benthics*), or in deep-water habitats (*deep benthics*). Finally, an ecogroup living over rocky habitats evolved (*mbuna*). The riverine generalist *Astatotilapia calliptera* (AC) living in the border regions of the lake and surrounding rivers, is the seventh ecogroup, and is thought to represent the ancestral lineage that seeded the three primary Lake Malawi lineages. Subsequent radiations within these ecogroups based on trophic specializations and sexual selection further diversified the species flock<sup>29,31</sup>.

Here, we investigate the role large inversions played in the Lake Malawi radiation using new large molecule technologies. We identify four single large inversions and two double inversions, ranging from 9.9Mb to 20.6Mb in size. A fifth single inversion, composed of just one of the two tandem inversions on chromosome 20, was also identified, indicating a serial set of rearrangements created this structural variant. These inversions primarily segregate by ecogroup, with no inversions found in the *mbuna* and AC samples, one inversion fixed in *Rhamphochromis*, three inversions fixed in *Diplotaxodon*, and all six segregating within the *benthic* sub-radiation. The evolutionary histories of these inversions are inconsistent with the species phylogeny and suggest that the three *Diplotaxodon* inversions spread into the *benthics* via hybridization around the time the benthic lineage diverged from the ancestral riverine species. The three additional inversions are found primarily in *benthic* species living in deep water habitats. We provide evidence that three of the inversions are involved in sex determination in a subset of *benthic* species. Our work provides a framework for understanding the role of inversions in the adaptive radiation of Lake Malawi cichlids, and suggests multiple roles in the establishment of ecogroups, adaptation to deep water habitats, and in controlling traits under sexual selection.

## RESULTS

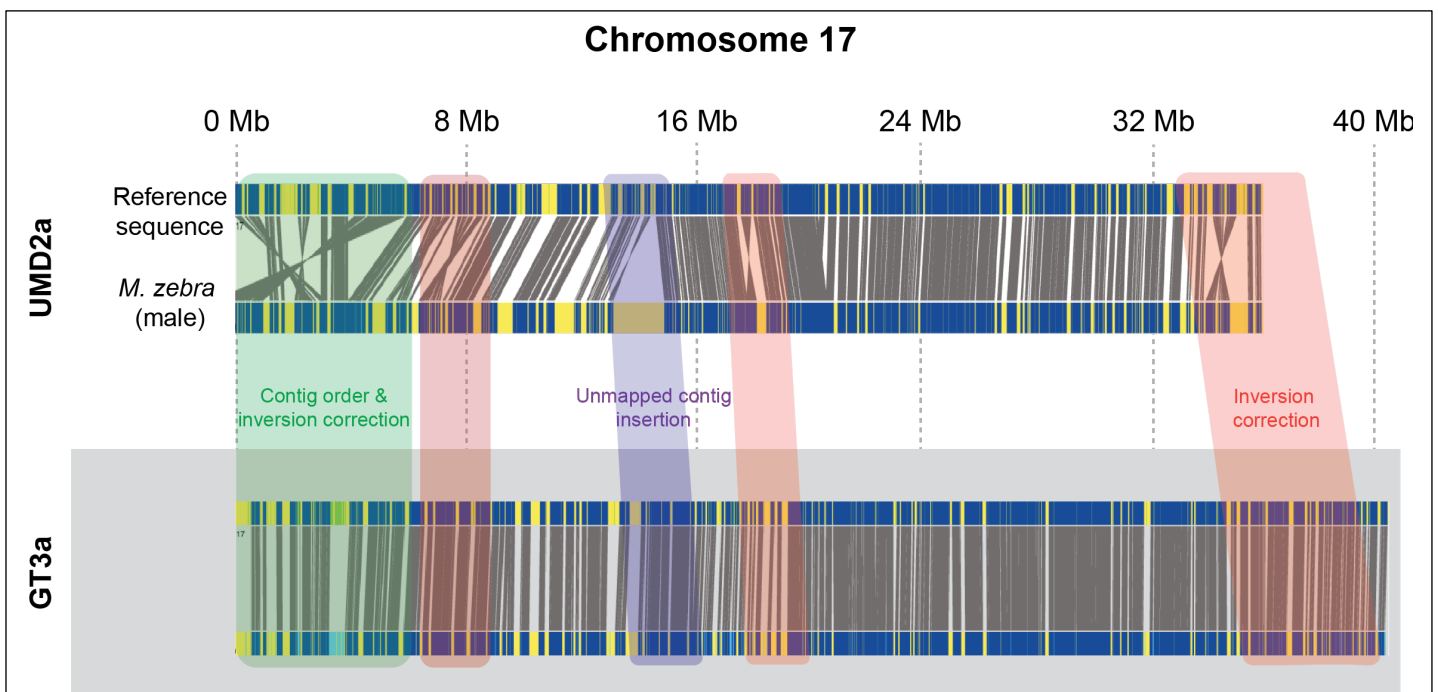
### *Improvement of the Metriaclima zebra reference*

The *mbuna*, *Metriaclima zebra*, was the first species sequenced as a reference genome for Lake Malawi cichlids<sup>38</sup>. The current version of its genome, called M\_zebra\_UMD2a, has been created by assembling PacBio reads into contigs which were anchored to 22 chromosomes using recombinant maps generated using crosses between different Lake Malawi species<sup>39</sup>. While the overall quality of this reference is high compared to other non-model organisms, a total of 20.5% of the DNA was not anchored to a chromosome. Additionally, the orientation of contigs can be ambiguous given how contigs were anchored to linkage groups using recombinant maps. These imperfections could complicate inversion discovery, so we first sought to improve the reference sequence.

The Bionano Saphyr system generates large (50 kb - 1 Mb) DNA molecules and labels a specific motif (CTTAAG) with a fluorescent probe<sup>40</sup>. While these molecules do not provide sequencing information, the distance between observed fluorescent loci can be used to assemble multiple molecules into large maps, with lengths spanning up to an entire chromosome. Comparison of these maps with an *in-silico* digest of the reference sequence identifies large structural changes in the observed sample or errors in the reference genome. We processed DNA collected from the blood of three *Metriaclima zebra* individuals

117 using the Bionano pipeline and assembled DNA molecules into optical maps. We identified many  
118 discrepancies between the observed optical maps and the M\_zebra\_UMD2a assembly (**Figure 1**). These  
119 differences fall into two primary categories: 1) contigs that were incorrectly ordered or oriented and 2)  
120 unplaced contigs that Bionano anchored to specific chromosomes. The first category is consistent with  
121 errors in the assembly caused by ambiguity in placing contigs into a linkage group using recombination  
122 maps. Our results indicate that 178 of the contigs are incorrectly oriented and 133 of the unanchored  
123 contigs could be placed onto specific chromosomes.

124 To further improve the reference, we used PacBio HiFi reads generated from an additional male *Metriaclima*  
125 *zebra* individual to create a new genome assembly, taking advantage of the low overall error rate (< 0.1%)  
126 of HiFi reads compared to previous generation PacBio technologies. Contigs generated from these data  
127 were combined with Bionano maps to generate a new hybrid genome assembly, which we refer to as the  
128 M\_zebra\_GT3a assembly. The overall quality of the assembly was excellent, with 933Mb of 962 Mb of DNA  
129 assigned to the 22 linkage groups, (**Table S1**) reducing the amount of unplaced contigs from 196 Mb to 28  
130 Mb. The N50 length was 32.542Mb, with many of the scaffolds containing entire chromosomes.

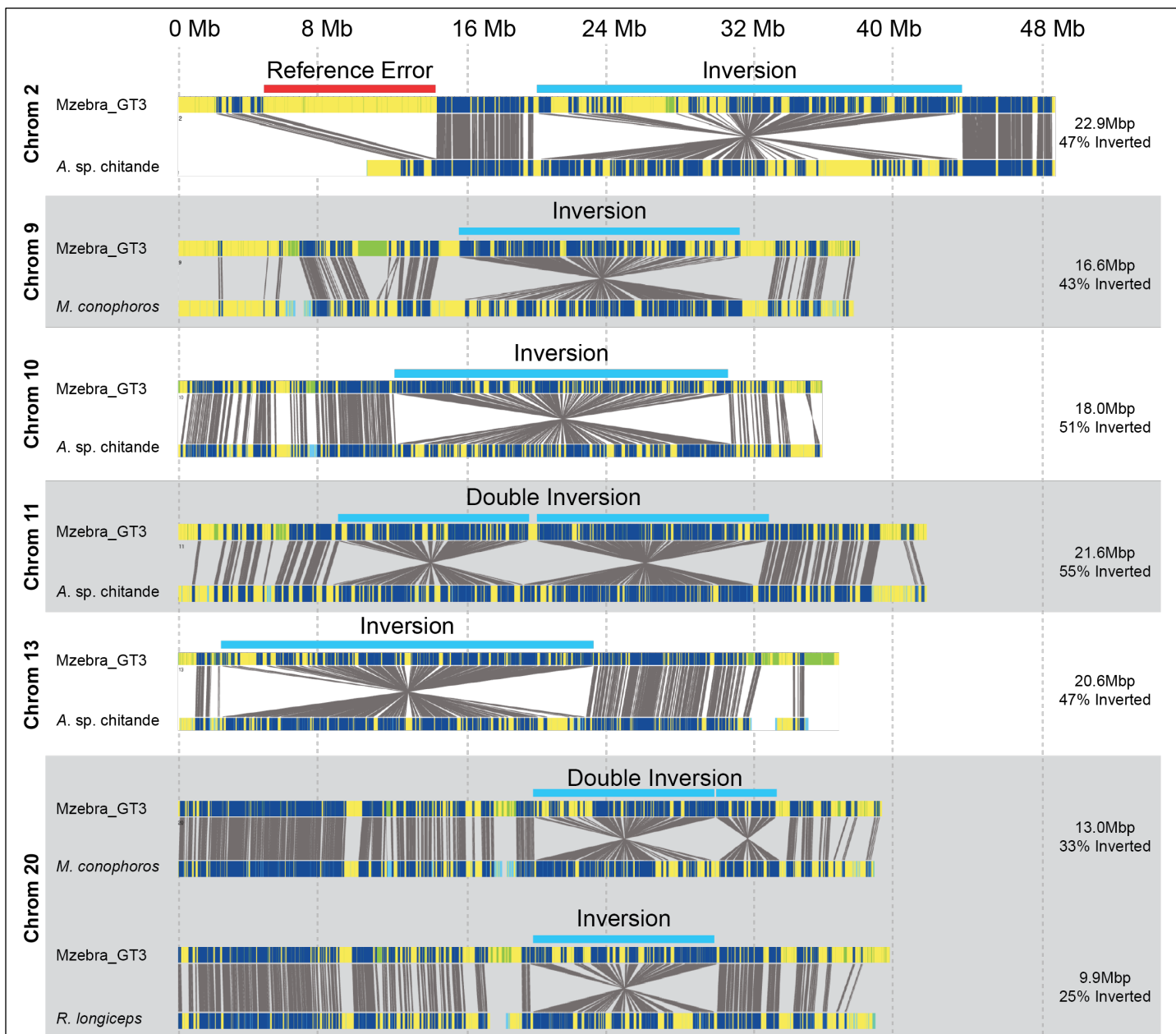


131 **Figure 1. Generation of a new high quality reference genome for *Metriaclima zebra*.** Example alignment of Bionano  
132 maps created from blood isolated from a single *Metriaclima zebra* male to the M\_zebra\_UMD2a reference (top) or the  
newly generated M\_zebra\_GT3a reference (bottom). The M\_zebra\_GT3a reference was created using a combination of  
PacBio HiFi reads and Bionano molecules. A summary of the improvements in the new genome is found in **Table S1**.

131

### 132 *Identification of six large inversions from 11 species of Lake Malawi cichlids*

133 Armed with the improved reference, we used the Bionano system to test eight different species of Lake  
134 Malawi cichlids that were currently used for experimental studies in the Streelman or McGrath laboratories  
135 (**Table 1 and Table S2**). We tested four *mbuna* species (*Pseudotropheus demasoni*, *Cynotilapia zebroides*,  
136 *Labeotropheus fuelleborni*, and *Labeotropheus trewavasae*), two *shallow benthic* species (*Mchenga*  
137 *conophoros* and *Nyassachromis prostoma* “Orange Cap”), a single *deep benthic* species (*Aulonocara* sp.  
138 ‘*chitande type north*’ Nkhata Bay), and a single *utaka* species (*Copidachromis virginalis*). To expand our  
139 coverage of the Lake Malawi radiation, we also obtained samples for a single species of the *Diplotaxodon*  
140 ecogroup, *Diplotaxodon limnothrissa*, and a single species of the *Rhamphochromis* ecogroup,  
141 *Rhamphochromis longiceps*. Finally, we also tested a single male *Protomelas taeniolatus shallow benthic*



**Figure 2. Identification of six large single or double inversions segregating within Lake Malawi cichlids.**

Representative alignments of inversions or double inversions identified from blood samples obtained from 30 individuals from eleven species. For each inversion, the top shows predicted motifs from an *in-situ* digest of the reference genome, the bottom shows motifs identified using Bionano molecules obtained from an individual of the species indicated on the left, and the grey lines indicate matching motifs based upon predicted and observed distances. Single inversions were identified on chromosomes 2, 9, 10, 13, and 20. Tandem double inversions were identified on chromosomes 11 and 20. The estimated length and percentage of the chromosome spanned by the inversion is shown on the right. The single inversion on 20 identified in *Rhamphochromis longiceps* has the same position as the first inversion of the double inversion on 20 identified from *Mchenga conophoros*. A list of all samples and their inversion genotype is in Table 1. Note that an error in the reference genome in chromosome 2 is indicated in the first panel.

individual based on its position in principal component analysis of SNVs described below. In total, we were able to obtain samples from 6 of the 7 ecogroups, excluding *Astatotilapia calliptera*.

From these 31 individuals, we identified six new inversions ranging from 9.9 - 22.9 Mb in size (**Figure 2**): a single inversion on chromosomes 2, 10, 11, and 13, and a double inversion on chromosomes 11 and 20. Interestingly, regarding the double inversion on chromosome 20, we also identified four *Rhamphochromis* individuals that carried only the first inversion, lacking the 4.1 Mb long second part of this rearrangement

(**Figure 2**). This indicates that the tandem inversions on 20 were formed by at least a two-part process, most parsimoniously with the left arm of the inversion (20a) forming first followed by the right inversion (20b) occurring in that genetic background.

The distribution within these species is summarized in **Table 1**. The most common structural rearrangements were the single inversion on chromosome 9, the double inversion on chromosome 11, and the double inversion on chromosome 20. These were found in a combination of *benthic* and *Diplotaxodon* individuals. The position of these three inversions were qualitatively consistent with regions of low recombination rate previously identified in an intercross between a *mbuna* and *benthic* individual, as expected, since inversions suppress recombination between inverted and non-inverted haplotypes<sup>22</sup>. The remaining three inversions (2, 10, and 13) were only found in the *Aulonocara* sp. 'chitande type north' *Nkhata Bay* individuals. While for the majority of samples the inversion genotype was homozygous, for five samples optical maps supported both inverted and non-inverted haplotypes, indicating the inversions were heterozygous in those individuals (**Figure S1**).

To provide further support for these structural rearrangements, we generated *de novo* hybrid genome assemblies of three *Aulonocara* sp. 'chitande type north' *Nkhata Bay* individuals, which carry all six of the identified inversions, and two additional *Metriaclima zebra* individuals using a combination of PacBio HiFi reads and Bionano maps. The overall quality of these assemblies is high, with large N50 values and a small number of individual contigs (**Table S3**). We performed whole genome alignments between the five new

**Table 1.** Inversion genotypes of 31 samples based on Bionano data.

Species	Ecogroup	#	Inversion						Notes
			2	9	10	11	13	20	
<i>Cynotilapia zebroides</i> 'Cobue'	<i>Mbuna</i>	1m/1f	-	-	-	-	-	-	
<i>Labeotropheus fuelleborni</i>	<i>Mbuna</i>	1m/1f	-	-	-	-	-	-	
<i>Labeotropheus trewavasae</i>	<i>Mbuna</i>	1m/1f	-	-	-	-	-	-	
<i>Metriaclima zebra</i>	<i>Mbuna</i>	3m/1f	-	-	-	-	-	-	Reference species
<i>Pseudotropheus demasoni</i>	<i>Mbuna</i>	1m/1f	-	-	-	-	-	-	
<i>Copidachromis virginalis</i>	<i>Utaka</i>	1m/2f	-	X	-	-*	-	X	One female was het on 11
<i>Protomelas taeniolatus</i>	<i>Shallow benthic</i>	1m	-	X*	-	-	-	X	The male was het on 9
<i>Mchenga conophoros</i>	<i>Shallow benthic</i>	1m/2f	-	X	-	X	-	X	
<i>Nyassochromis prostoma</i> 'Orange Cap'	<i>Shallow benthic</i>	1m/1f	-	X	-	X*	-	X	The male was het on 11
<i>Aulonocara</i> sp. 'chitande type north' <i>Nkhata Bay</i>	<i>Deep benthic</i>	3m/1f	X	X	X*	X	X	X	All three males were het on 10
<i>Diplotaxodon limnothrissa</i>	<i>Diplotaxodon</i>	2m/0f	-	X	-	X	-	X	
<i>Rhamphochromis longiceps</i>	<i>Rhamphochromis</i>	3m/1f	-	-	-	-	-	X*	All samples had the single left inversion on 20

167 assemblies and the GT3a reference genome (**Figure S2-S7**). These alignments show that the new  
168 assemblies support the presence of all six inversions in the *Aulonocara sp. 'chitande type north' Nkhata*  
169 *Bay* individuals. For most inversions, a single scaffold covered the entire inversion with breakpoints in  
170 locations consistent with the Bionano maps. The two new control assemblies were consistent with the  
171 M\_zebra\_GT3a reference and did not show any evidence of inversions in those regions. In other regions of  
172 the genome, however, there were differences in alignment of the M\_zebra\_GT3a assembly and the  
173 *Metriaclima zebra* genomes, most prominently in the left arm of chromosome 2 (**Figure 2**) and in the highly  
174 repetitive chromosome 3, which likely indicate some errors in the M\_zebra\_GT3a assembly.

175 We refined the location of the breakpoints for each inversion using the genomic alignments. However, the  
176 presence of a large amount of repetitive DNA in regions near the breakpoints prevented us from defining  
177 the inversion boundaries at single base pair resolution.

178 We also took advantage of a genome assembly for a *Rhamphochromis sp 'chilingali'* individual  
179 (fRhaChi2.1) created by the Darwin Tree of Life Project produced using PacBio data and Arima2 Hi-C data  
180 (Accession PRJEB72870). The genome alignment between the fRhaChi2.1 assembly and the GT3a  
181 assembly confirmed a single inversion on chromosome 20, consistent with the Bionano data (**Figure S2-**  
182 **S7**).

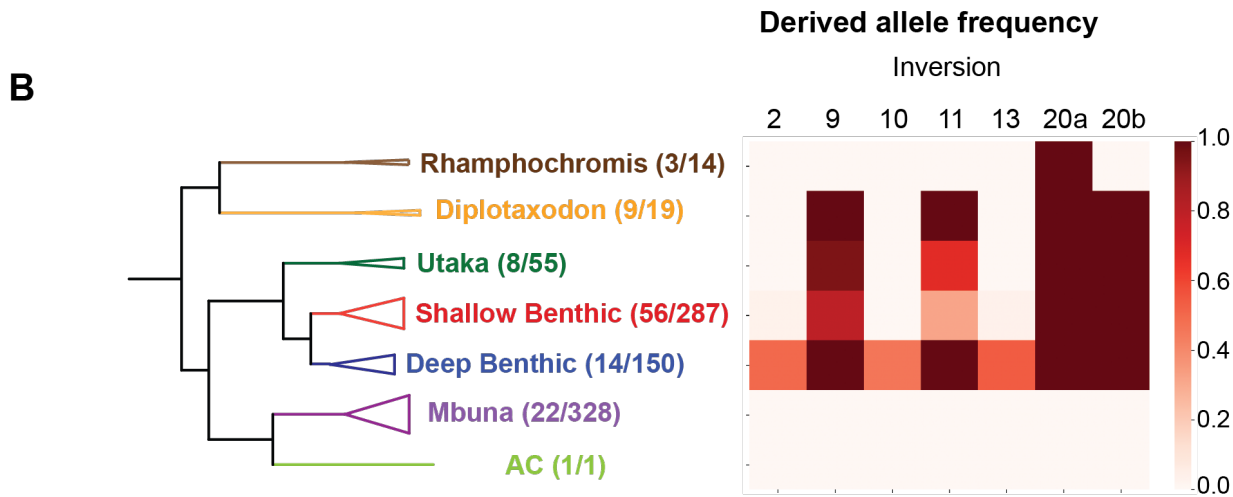
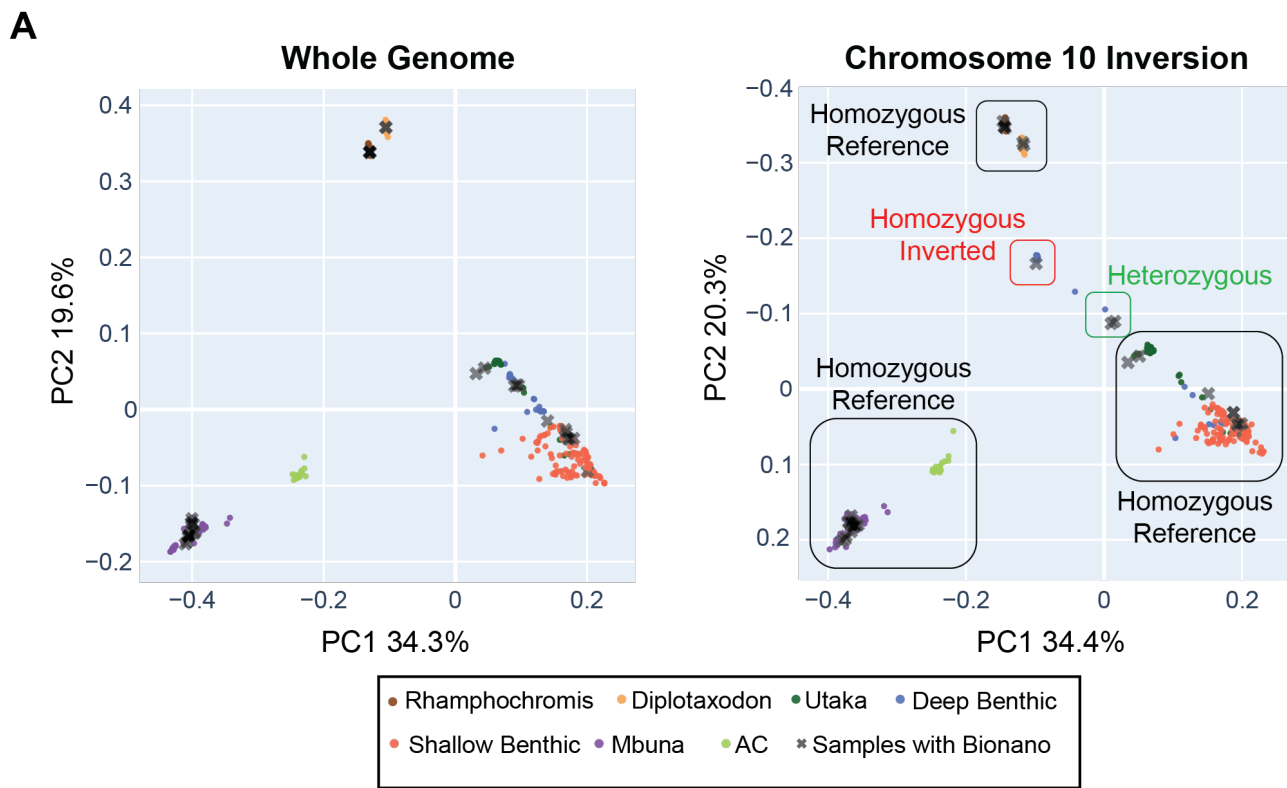
183 For each inversion, we determined the ancestral haplotype using two species as outgroups, *Oreochromis*  
184 *niloticus* and *Pundamilia nyererei*. *Oreochromis niloticus*, or Nile tilapia, is an important food fish that's  
185 estimated to have diverged from Lake Malawi cichlids 14.1 to 30 million years ago while *Pundamilia*  
186 *nyererei* is a more closely related haplochromine found in Lake Victoria<sup>38,41</sup>. For this analysis we used a  
187 published whole genome assembly for *Oreochromis niloticus* (UMD\_NMBU)<sup>39</sup> and generated a hybrid  
188 genome assembly for *Pundamilia nyererei* using a previously published whole genome assembly<sup>38</sup>  
189 combined with Bionano maps we produced from a single individual (**Figure S8**). For all six inversions, the  
190 reference haplotype found in the GT3a assembly represents the ancestral state.

#### 191 *Inference of the segregation of these inversions in Lake Malawi using short-read sequencing*

192 The distribution of these inversions within the lake can provide important information to their function.  
193 Genotyping inversions using short read sequencing is possible because each haplotype is on an  
194 independent evolutionary trajectory due to the suppression of recombination. Population genetics  
195 approaches take advantage of the fact that inversions create patterns of inheritance between SNVs that  
196 are readily detected by approaches such as Principal Component Analysis (PCA)<sup>42</sup>.

197 We analyzed short read sequencing data from the 31 samples with Bionano data along with 297 wild  
198 individuals published by the Durbin lab<sup>29</sup>, 15 wild individuals published by the Streelman lab<sup>30,43</sup>, and 22  
199 wild individuals sequenced for this paper, aligning reads and calling small variants against our new  
200 M\_zebra\_GT3a reference genome. PCA analysis was used to analyze the major axes of variation across the  
201 whole genome and within the six inversions. We separately analyzed the left and right inversions on 11 and  
202 20 (the two inversions on 11 were identical in their distribution). For the whole genome, the distribution of  
203 ecogroups within the first two principal components matched the previously reported results of Malinsky  
204 et al<sup>29</sup>. Individuals from four of the ecogroups - *Rhamphochromis*, *Diplotaxodon*, *Astatotilapia calliptera*,  
205 and *mbuna* - formed four distinct clusters (**Figure 3A** and **Figure S9**). A fifth, more diffuse cluster,  
206 composed of *shallow benthics*, *deep benthics*, and *utakas*, was also present. The broader distribution of  
207 the *benthic* ecogroups in PC space is consistent with both their recent separation as well as large amounts  
208 of gene flow thought to have occurred between these different species<sup>29,35</sup>.

209 The PCA plots using the SNVs within the six inversions often violated the pattern of the whole genome PCA  
210 (**Figure 3A** and **Figure S10-S17**). Using samples that were sequenced with both Illumina and Bionano



**Figure 3. Distribution of inversions within the Lake Malawi ecogroups.** Principal component analysis was used to analyze SNVs identified using whole genome sequencing from 365 samples to genotype the inversion haplotype for all six inversions. **(A)** We illustrate the approach for genotyping the chromosome 10 inversions. The PCA plot for the entire genome is shown on the left, with each sample colored by its ecogroup. Individuals with Bionano data are shown as grey Xs. The deep benthic/shallow benthic and utaka individuals that cluster together for the whole genome analysis split into multiple clusters when analyzed using the SNVs that fall within the chromosome 10 inversion, and these clusters were assigned to inversion genotypes using samples that were genotyped using Bionano data. To make comparisons of the whole genome PCA and chromosome 10 PCA plots easier, we reversed the y-axis on the right panel. The genotypes of each sample can be found in Supplemental Table 2. Interactive PCA plots for each inversion are included in **Figures S9-S17**. **(B)** The derived allele frequency was calculated for each inversion within the seven ecogroups. The number of species that were genotyped and the estimated number of species that live in Lake Malawi are listed in the whole-genome phylogeny (left).

technologies, we assigned each cluster to specific inversion genotypes (**Table S2**). We did not have Bionano data for two clusters (one on 9 and one on 11 – see **Figures S11, S13, and S14**), so we



213 supplemented our dataset by adding Bionano data from an additional *shallow benthic* individual  
214 (*Protomelas taeniolatus*) that fell into both clusters. The PCA clusters varied in their complexities for each  
215 inversion depending largely on their distribution within the *benthic* ecogroup (more on this below) and  
216 included both homozygous genotypes and heterozygous genotypes. We summarize the overall distribution  
217 of the inversions within each of the ecogroups in **Figure 3B**.

218 All specimens of the *mbuna* and *AC* ecogroup were fixed for the homozygous non-inverted alleles for all six  
219 inversions. In *Rhamphochromis*, all the inversions were homozygous non-inverted except for the left arm  
220 of 20 which was fixed for the single inversion state. The *Diplotaxodon* individuals all carried the non-  
221 inverted alleles for the inversions on 2, 10 and 13 and were fixed for the inversions on 9, 11, and 20.

222 The distribution of these inversions was the most complicated within the benthics. The inversions on 2, 10,  
223 and 13, which were identified in *Aulonocara* sp. ‘chitande type north’ Nkhata Bay, were mostly found  
224 together in other *deep benthics*, including the *Alticorpus* (*geoffreyi*, *macrocleithrum*, *peterdaviesi*),  
225 *Aulonocara* (*blue chilumba*, *gold* and *minutus*), and *Lethrinops* (*gossei*, *longimanus* ‘redhead’, and sp.  
226 *olivera*) genera. The inversions on 2 and 13 were also found in the shallow benthic species *Placidochromis*  
227 *longimanus*. However, not all the *deep benthics* carried the three inversions, including five *Aulonocara*  
228 species (*baenschi*, *getrudae*, *steveni*, *stuartgranti*, and *stuartgranti* Maisoni) and a single *Lethrinops*  
229 species (*longipinnis*). The distribution of these three inversions in *benthic* species known to inhabit the  
230 deepest habitats are suggestive for a role in depth adaptation.

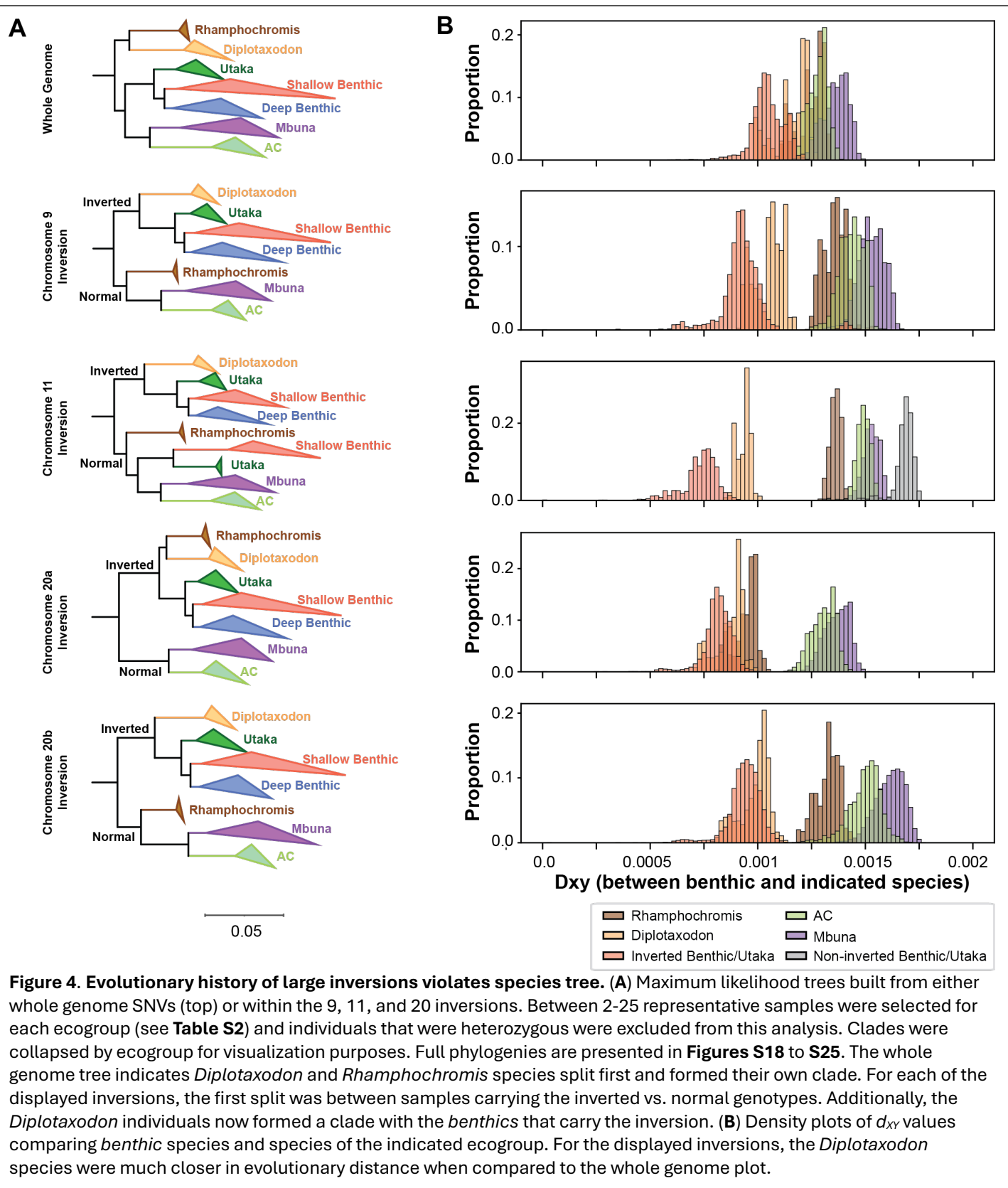
231 The double inversion on 20 was homozygous inverted in all benthic individuals. For the single inversion on  
232 9, many *benthics* were homozygous for the inverted haplotype. However, an additional PC cluster was  
233 present, composed of 64 *shallow benthic* and *utaka* individuals heterozygous for the non-inverted  
234 haplotype (**Figure S1**). Interestingly, no *benthic* individuals were identified as homozygous for the non-  
235 inverted state.

236 Finally, the distribution of the chromosome 11 double inversion was also determined. Individuals of all  
237 three genotypes (homozygous non-inverted, homozygous inverted, and heterozygous) were observed  
238 (**Figure 3a**).

#### 239 *Evolutionary history of 9, 11, and 20 suggest introgression from Diplotaxodon into benthic ancestors*

240 The distribution of these inversions conflicts with the species tree, suggesting that these inversions could  
241 have spread via hybridization. To determine their evolutionary history, we built phylogenies for the whole  
242 genome as well as for each inversion using SNVs that fell within each of the six inversions using  
243 representative species from each ecogroup (73 total – **Table S2**) as well as two previously sequenced  
244 *Pundamilia nyererei* individuals as an outgroup (**Figure 4A** and **Figures S18 – S25**)<sup>44,45</sup>. To avoid issues with  
245 haplotype phasing, we focused on individuals that were homozygous for one of the two inverted genotypes.

246 While the whole genome phylogeny matched the expected topology based upon our understanding of the  
247 Lake Malawi species tree, the topologies of each of the inversions showed key differences between the  
248 separation of ecogroups. The *benthic* group often split into separate clades, correlated with their inversion  
249 genotype. Prominently, we found the three inversions on 9, 11, and 20 showed key differences in the  
250 relationship between the *Diplotaxodon* and *benthic* ecogroups. While *Diplotaxodon* is a sister group to  
251 *Rhamphochromis* at the species level, for these three inversions, the *Diplotaxodon* individuals formed a  
252 clade with the *benthics* that also carried the inversion. Additionally, genetic distance between  
253 *Diplotaxodon* and *benthics* carrying the inversion was also much smaller than the rest of the genome  
254 (**Figure 4B**). Both the phylogenies and genetic distances between ecogroups are consistent with the  
255 introgression of the 9, 11, and 20 inversions from the *Diplotaxodon* to *benthic* ancestors at the time of the  
256 benthic radiation.

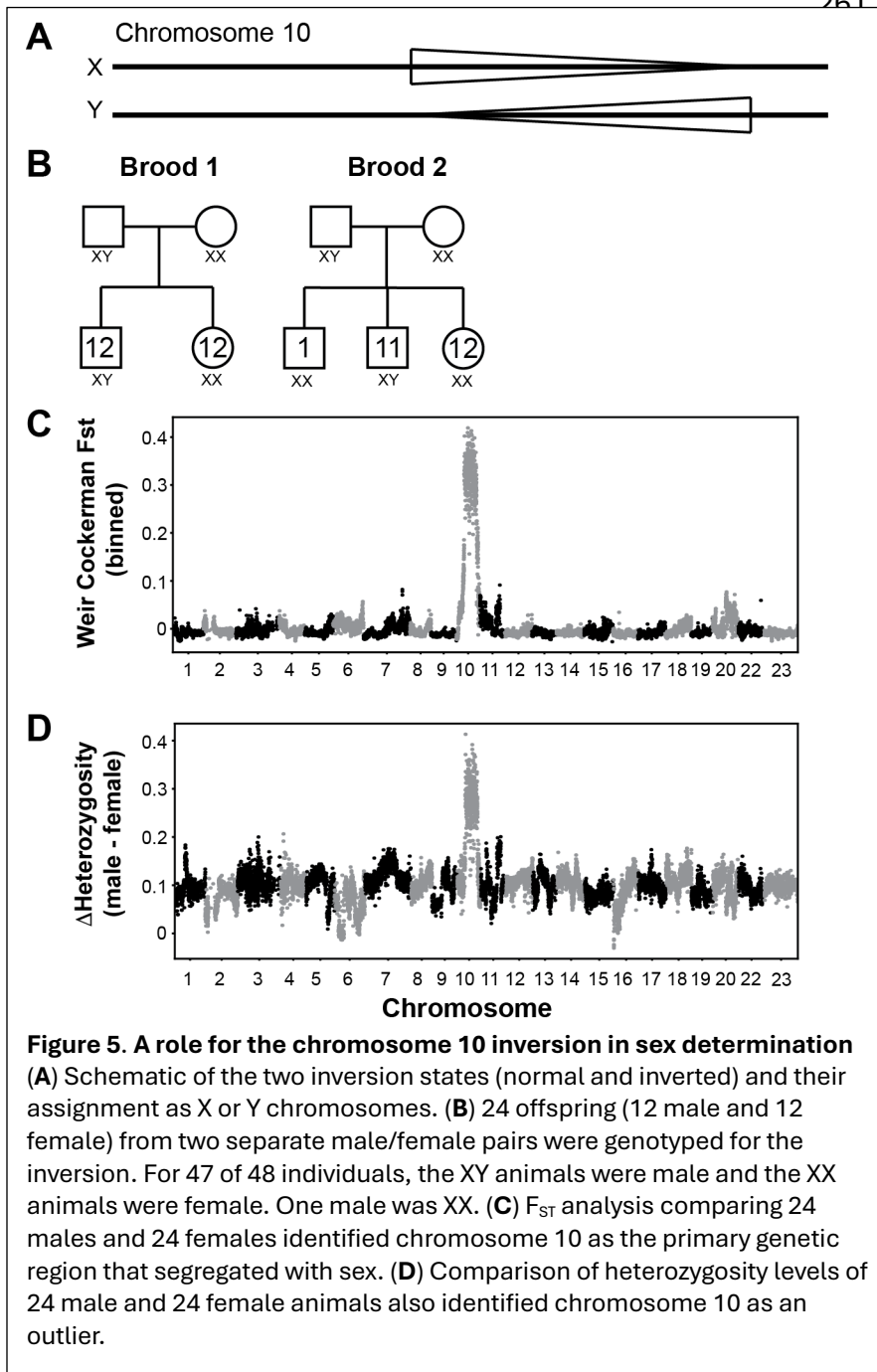


257  
258

#### A role for large inversions in sex determination within the benthic ecogroup

259  
260

The evolution of sex determination is often associated with the presence of inversions to repress recombination between the sex chromosomes<sup>46,47</sup>. We tested whether these inversions could play a role



sex. Males also had a much higher heterozygosity in this region than females (**Figure 5D**). Altogether, these results indicate that the inversion on 10 acts as an XY sex determiner in the *Aulonocara sp. 'chitande type north' Nkhata Bay* species.

We also tested eight laboratory *Nyassachromis prostoma* “Orange Cap” individuals, as our Bionano data identified a male that was heterozygous for the inversion on 11 and a female that was homozygous for the inversion. We performed short-read sequencing on three males and five females growing in the lab, genotyping the inversion using PCA analysis. All three males were heterozygous for the inversion while the five females were homozygous, a significant association between genotype and sex (Fisher exact test p-value = 0.0179). This suggests a separate XY system segregates in *Nyassachromis prostoma* “Orange Cap” using the inversion on chromosome 11.

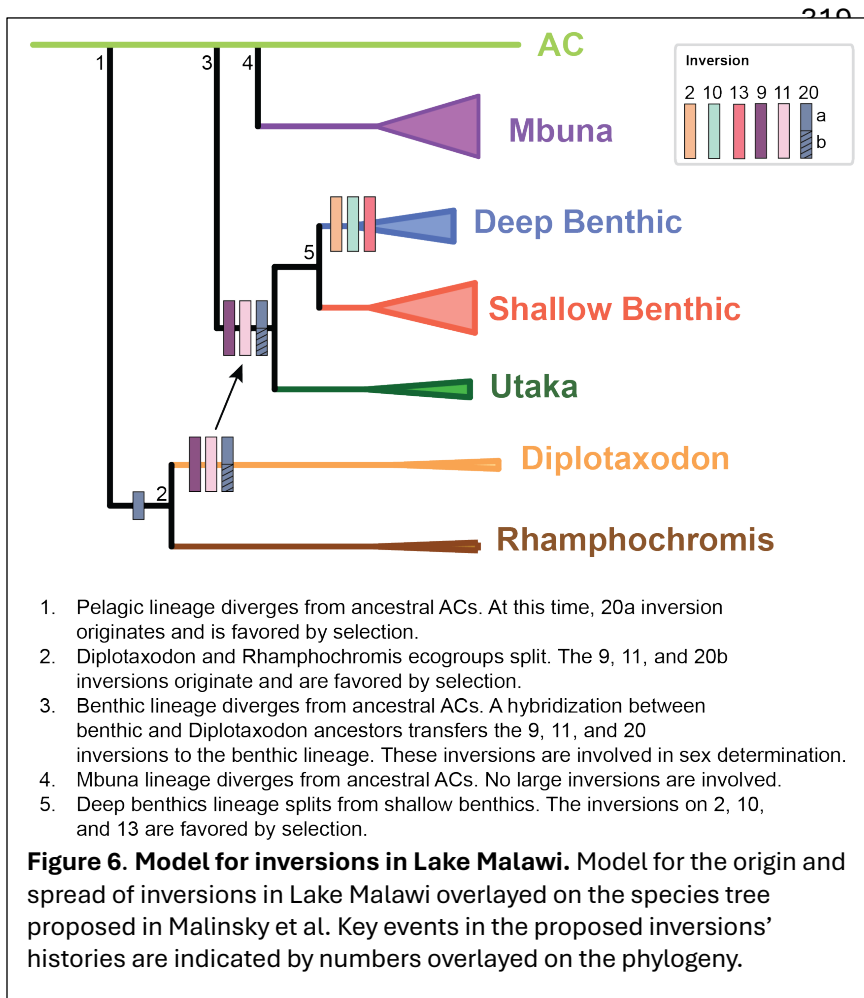
in sex determination in *benthic* species, as we identified individuals that were heterozygous for some of these inversions (**Table 1**). We started with the inversion on 10, as our Bionano data identified three male *Aulonocara sp. 'chitande type north' Nkhata Bay* individuals that were heterozygous and a single female that was homozygous for the inverted haplotype. To test whether this inversion acted as an XY system in this species, we tested two separate broods (Brood 1; 2 parents, 24 offspring and Brood 2: 2 parents, 24 offspring) currently growing in the laboratory, sequencing 12 males and 12 females from each brood with Illumina paired-end sequencing. We used PCA to genotype each of the 52 animals (4 parents and 48 offspring) for the chromosome 10 inversion genotype. The association in the offspring between the inversion 10 genotype and sex was almost perfect: 24 females and 1 male were homozygous for the inversion, while 23 males were heterozygous for the inversion (**Figure 5B**).

We also used binned Weir Cockerham  $F_{ST}$  analysis, grouping male and female individuals to scan the genome for regions that were associated with sex (**Figure 5C**). The inversion on 10 was the region of the genome that was in strongest linkage disequilibrium with

Finally, for the inversion on nine, there was a cluster of heterozygous individuals that included 7 individuals that we had sex information for. All 7 of these individuals were male (**Table S2**). While this association is not significant, it is also suggestive that the inversion on 9 could play a role in sex determination as well.

## DISCUSSION

Here, we characterize the presence of large inversions found in the Lake Malawi cichlid flock, identifying four single inversions and two double inversions, providing a framework to address the role structural variants played in this lake's adaptive radiation (**Figure 6**). Inversions can have multiple roles in the evolutionary process: 1) they can capture beneficial alleles responsible for adaptation to local habitats, 2) they can capture alleles that create pre- and post-zygotic barriers that prevent gene flow during speciation, and 3) they can capture alleles that influence sex-specific phenotypes that are linked to sex-determination systems<sup>1,3,4</sup>. Future work is now possible to determine how these inversions contribute to the evolution of Lake Malawi cichlids through these mechanisms.



The overall size of these inversions is large, together capturing >10% of the genome (~112 Mb). Apart from the *benthics*, the inversions are fixed in each ecogroup in the inverted or non-inverted state. Such a distribution would be expected in the local adaptation hypothesis if they were involved in the adaptation of these ecogroups to their respective lake habitats. In this scenario, the inversions would be under selection prior to speciation, and would play roles in: 1) the split of the pelagic group from the ancestral riverine lineage (the left arm of the inversion on 20), 2) the diversification of the *Rhamphochromis* and *Diplotaxodon* groups (the inversions on 9, 11, and right arm of 20), and 3) the adaptation of *benthics* to deeper water habitats (the inversions on 2, 10, and 13). The role of inversions in local adaptation between spatially separated populations has been well characterized<sup>6,8</sup>. If two populations are undergoing adaptive divergence in the presence of migration

and gene flow, theory indicates that inversions will be selected for when they capture and link adaptive alleles together<sup>48</sup>. Given gene flow was likely common during the formation of the ecogroups<sup>49</sup>, we expect inversions that captured adaptive alleles that increase fitness in a habitat-specific context to feature prominently in ecogroup inception. This hypothesis predicts that adaptive alleles for each ecogroup are enriched within the inversions, which can be tested in Lake Malawi cichlids, as different species can interbreed to test whether the inversions are associated with differences in phenotypes.

While the role of the inversions in most of the ecogroups could potentially be explained by local adaptation, their role in the benthic sub-radiation appears to be more complicated. Inverted and non-inverted

haplotypes for five of the six inversions are segregating within this group. Additionally, the presence of the 9, 11, and 20 inversions within *benthics* violates the whole species phylogeny and is consistent with its spread from the *Diplotaxodon* to a very early *benthic* ancestor via introgression between ancestors of these two groups. Hybridization has been proposed to fuel adaptive radiations of cichlids in all three great African lakes (Malawi, Tanganyika, and Victoria) and our work suggests that an ancient hybridization event between early *benthic* and *Diplotaxodon* ancestors occurred prior to the benthic radiation resulting in the 3 inversions (as well as other regions of the genome) spreading into what ultimately became the *benthic* ecogroup<sup>26,27,36</sup>. It will be interesting to study the role of these three inversions in *benthic* species. *Diplotaxodon* are pelagic, deep-water dwellers that feed on plankton or small fish while *benthics* live near the shore over sandy and muddy habitats<sup>29</sup>. It is not obvious how the traits of the *Diplotaxodon* would benefit *benthic* animals that live in such dissimilar environments. However, the initial habitats that the *benthic* ancestors utilized is not known and potentially their initial habitat was in deeper waters than the AC ancestors. Whether or not these inversions were involved in the benthic radiations, they are not strictly required for large radiations as the *mbuna* group, with an estimated >300 species, carries zero large inversions<sup>29</sup>. This could reflect the relative importance of allopatric vs. sympatric speciation in the two radiations, as the rocky habitats the *mbuna* inhabit are often spatially separated throughout the lake.

The inversions we identified also play a role in sex determination in the *benthics*, with at least three of the inversions (9, 10, and 11) acting as XY systems to control sex determination. Despite the importance of the maintenance of distinct, reproductively compatible sexes, sex determination systems can be evolutionarily labile, and Lake Malawi cichlids show exceptional variability, with at least 5 different genetic sex determination systems described to date<sup>50,51</sup>. While we implicated multiple inversions in sex determination, our work also underscores how dynamic sex determination is within the *benthics*. While many benthic species carry the different inversions, their genotypic distribution in most species is inconsistent with a role for sex determination. Effectively, the genetic determination of sex determination largely must be determined in a species by species way, and individual species may be segregating multiple sex determination systems. Additionally, these inversions could have played a role in sex determination in the ancestors of pelagic ecogroups, which could have contributed to their spread before sex chromosome turnover ended balancing selection on the inverted/non-inverted haplotypes and subsequent fixation of these inversions. Identification of the causal alleles for sex determination on these inversions will allow us to resolve this question.

The presence of non-inverted haplotypes in the *benthics* could have been retained from the time the lineage initially formed, with both haplotypes under balancing selection if they were used for sex determination. A recent preprint, however, which identified 5 of the 6 inversions presented here via short-read sequencing of 1375 individuals sampled from Lake Malawi, presented an alternative model<sup>52</sup>. They proposed the non-inverted haplotypes of the 9 and 10 inversions were introgressed back into the *benthics* from riverine species within and outside of the lake. Similarly, the non-inverted haplotypes on chromosomes 2, 10, and 13 in the *deep benthics* were introgressed from *shallow benthics*. In this model, it is unclear what evolutionary forces drove these back introgressions of the non-inverted alleles and when they developed a role in sex determination.

The selective forces on sex chromosomes are strong, including balancing selection for sex ratio and sexually antagonistic selection on sexually dimorphic phenotypes. Lake Malawi cichlids show amazing sex differences, characterized by diversification in numerous phenotypes, including pigmentation, behavior, and morphology thought to contribute to species barriers. The use of inversions in sex-determination allows for the accumulation of alleles in linkage with the sex-determiner gene which will have sex-specific effects in the male individuals that carry the Y chromosomes. Interestingly in all three cases, the non-inverted haplotype acts as the Y, suggesting that different haplotypes could play independent roles in

399 evolution. For example, the inversion X haplotype on 10 could carry alleles responsible for adaptation to  
400 deep water habits, which would be carried by both sexes, while the non-inverted Y haplotype carries alleles  
401 under sexual selection.

402 This report provides an initial framework for understanding the role of inversions in the adaptive radiation  
403 of Lake Malawi cichlids. Future work identifying and characterizing causal alleles captured within these  
404 inversions will enable us to discern the specific mechanisms by which they contribute to speciation using  
405 this remarkable evolutionary system as a model for the evolution of all animals.

## 406 ACKNOWLEDGEMENTS

407 We thank Krish Roy for the acquisition of the Bionano Saphyr system, Manasi Pimpley for assistance in  
408 adapting existing Bionano kits to cichlid tissue, and the Molecular Evolution Core Laboratory at the Parker  
409 H. Petit Institute for Bioengineering and Bioscience at the Georgia Institute of Technology for the use of  
410 their shared equipment, services, and expertise. This work was supported in part by NIH R35 GM139594  
411 and the Nelson and Bennie Abell Professorship to P.T.M., R01GM144560 to J.T.S., and U.S. National Science  
412 Foundation (DEB-1830753) to T.D.K..

## 413 METHODS

414 **Samples.** Wild-caught cichlids were acquired from Old World Exotic Fish Inc in 2022 and Cichlidenstadel  
415 in 2024. Lab-reared species were purchased from Southeast Cichlids, Old World Exotic Fish Inc., and  
416 Cichlidenstadel at various points in the past two decades (**Table S2**). Fish were delivered live to the Georgia  
417 Institute of Technology cichlid aquaculture facilities. All samples were collected following anesthetization  
418 with tricaine according to procedures approved by the Institutional Animal Care and Use Committee  
419 (IACUC protocol numbers A100029 and A100569).

420 Caudal fins were collected from subjects and flash frozen on powdered dry ice before storage at -80°C.  
421 Care was taken to minimize freeze-thaw cycles prior to DNA extraction via Bionano or Qiagen protocols. If  
422 necessary, fins were sectioned on a sterile aluminum block set in dry ice prior to digestion.

423 Blood was collected following rapid decapitation of anesthetized subjects. Wide-bore, low-retention  
424 pipette tips, 0.5M EDTA (pH 8.0, Invitrogen Cat# 15575-038), and Bionano Cell Buffer (Part Number 20374)  
425 prevented blood from clotting at room temperature. Samples were used immediately for Bionano DNA  
426 extraction (see below).

427 **Bionano.** Bionano DNA extractions were predominantly performed on fresh whole blood samples (see  
428 above). Blood concentration was determined via a manual logical count performed on a hemocytometer  
429 and 2 million cells were carried forward for DNA extraction with the Bionano SP-G2 Blood & Cell Culture  
430 DNA Isolation Kit (Part Number 80060) as per the Bionano Prep SP-G2 Fresh Human Blood DNA Isolation  
431 Protocol (Document Number CG-00005, Rev C). Because nucleated cichlid blood had high  
432 concentrations, steps 3 and 4 were skipped and samples were resuspended in Bionano Cell Buffer to a  
433 total volume of 200µL. Following extraction, samples were allowed to homogenize undisturbed at room  
434 temperature for at minimum 72 hours prior to quantification using the Qubit dsDNA Quantitation, Broad  
435 Range kit (ThermoFisher, Cat # Q32853).

436 Bionano data for a single *Protomelas taeniolatus* sample (PT\_2003\_m, see **Table S2**) was generated from  
437 fin tissue using a combination of the Bionano Prep SP Tissue and Tumor DNA Isolation Protocol  
438 (Document Number 30339, Rev A) and the Bionano Prep SP-G2 Fresh Human Blood DNA Isolation  
439 Protocol.

440 750ng of homogenized DNA was fluorescently labeled using the Bionano Direct Label and Stain-G2 (DLS-  
441 G2) Kit (Part Number 80046, Protocol - CG-30553-1, Rev E). Labeled DNA was quantified using the Qubit  
442 dsDNA Quantitation, High Sensitivity kit (ThermoFisher, Cat # Q32854) and loaded onto a flow cell within

443 a Bionano Saphyr Chip G3.3 (Part Number 20440). Samples were imaged on the Bionano Saphyr with at  
444 least 500Gbp of data collected per sample.

445 Molecules for each sample were assembled using the Bionano *de novo* assembly pipeline on Bionano  
446 Access Version 1.7. The preassembly setting was turned on, and the variant annotation parameters were  
447 deselected for every assembly. *De novo* maps were aligned to the Mzebra\_GT3a genome.

448 **PacBio HiFi sequencing and genome assembly.** PacBio *de novo* assemblies were generated from HiFi  
449 reads across three sequencing runs. The MZ\_GT3.3 and YH\_GT1.3 were sequenced following DNA  
450 extraction from frozen caudal fin tissue using the Qiagen MagAttract HMW DNA Kit (Cat. No. 67563). In the  
451 first sequencing run a *Metriaclima zebra* female was sequenced on a PacBio Sequel II system  
452 by the Georgia Genomics and Bioinformatics core. In the second run, the same the *Metriaclima zebra*  
453 female and an *Aulonocara sp. 'chitande type north' Nkhata Bay* female was sequenced on a PacBio Revio  
454 instrument by the HudsonAlpha Institute for Biotechnology. Reads from both runs were used to assemble  
455 the MZ\_GT3.3 genome and reads from the second run alone were used to assemble the YH\_GT1.3 genome  
456 (see details below).

457 Mzebra\_GT3a, MZ\_GT3.2, YH\_GT1.1, and YH\_GT1.2 were assembled using HiFi reads generated with DNA  
458 extracted from fresh whole blood and heart tissue. DNA was extracted via the PacBio Nanobind Tissue Kit  
459 RT (Part Number 102-208-000) following a modified DNA from animal tissue using the Nanobind®  
460 PanDNA kit protocol (Part Number 102-574-600). Blood and heart tissue were combined and  
461 homogenized together using a Qiagen TissueRuptor (Qiagen 9002755) and centrifugation g-force was  
462 halved during the DNA extraction. DNA fragments <25kb were removed using the PacBio Short Read  
463 Elimination kit (Part Number 102-208-300). Library preparation and sequencing was performed by the  
464 University of Maryland Institute for Genome Sciences on a PacBio Revio instrument.

465 HiFi reads from all samples were assembled using the Mabs *de novo* assembler (v2.28)<sup>53</sup> using the `mabs-  
466 hifiasm` algorithm and default parameters. Each assembly from `mabs-hifiasm` was evaluated using  
467 `Inspector (inspector.py v1.0.1)`<sup>54</sup> run on default parameters. The assemblies were error corrected with  
468 `inspector-correct.py (v1.0)` to resolve large structural errors. These error-corrected contigs from  
469 were uploaded to Bionano Access.

470 The MZ\_GT3.2, MZ\_GT3.3, YH\_GT1.1, YH\_GT1.2 and YH\_GT1.3 scaffold-level assemblies were generated  
471 by bridging gaps in their contigs using Bionano maps via the single-enzyme Bionano Hybrid Scaffold  
472 pipeline using default parameters. A unique set of Bionano maps was used to scaffold the contigs from  
473 each *de novo* assembly. The resulting hybrid scaffold `NCBI.fasta` file (which denotes gaps filled by  
474 Bionano maps with stretches of N nucleotides) was concatenated with any unscaffolded contigs from  
475 the error-corrected mabs assembly. These five scaffold-level assemblies have been deposited at  
476 DDBJ/ENA/GenBank (accessions TBD).

477 A scaffold-level genome was identically generated for Mzebra\_GT3a. These scaffolds were further aligned  
478 to the corrected UMD2a genome (see below) and anchored to linkage groups using D-GENIES<sup>55</sup> with the  
479 Minimap2 v2.28 aligner and “Many repeats” options. An anchored genome was output using the Query  
480 assembled as reference export option.

481 Note that the mitochondrial DNA sequence in Mzebra\_GT3a is the same as the mitochondria assembled  
482 in UMD2a. This Whole Genome Shotgun project has been deposited at DDBJ/ENA/GenBank under the  
483 accession JBEVYI000000000. The version described in this paper is version JBEVYI010000000.

484 **Correcting errors in M\_zebra\_UMD2a.** We generated Bionano molecules for 3 male and 1 female  
485 *Metriaclima zebra* subjects, assembled Bionano maps, and aligned them to the UMD2a reference genome.  
486 We reordered or reoriented contigs that were revealed as misassemblies or inversions, respectively, in all

four Bionano maps and also corresponded to breakpoints between contigs in UMD2a. By filtering interchromosomal translocation calls in the Bionano Access software, we were able to insert unmapped contigs into the established linkage groups.

**Illumina extractions and sequencing.** DNA was extracted from fresh or frozen fin tissues for short-read sequencing using the Qiagen MagAttract HMW DNA Kit (Cat. No. 67563) or the Qiagen DNeasy Blood & Tissue Kit (Cat No. 69504) using the Manual Purification of High-Molecular Weight Genomic DNA from Fresh or Frozen Tissue and Purification of Total DNA from Animal Tissues (Spin-Column) protocols respectively. DNA was delivered to the Georgia Tech Molecular Evolution Core where libraries were prepared with the NEBNext® Ultra™ II FS DNA Library Prep Kit for Illumina (NEB #E6177) using the Protocol for FS DNA Library Prep Kit (E7805, E6177) with Inputs  $\geq 100$  ng. A small subset of samples was sent to an external collaborator by the Molecular Evolution Core where library preparation was performed using the KAPA HyperPrep Kit (Roche, Material Number: 07962363001). Samples across all runs were sequenced on a NovaSeq 6000 instrument using v1.5 chemistry.

**Variant calling and PCA analysis.** Fastq files from Illumina sequencing were converted to UBAM format using the `gatk56 FastqToSam` algorithm. UBAM files were used for alignment to `Mzebra_GT3` using `bwa57 mem`; the `-M` and `-p` flags were used. To carry forward metadata from the unaligned BAMs, `gatk MergeBamAlignment` was called on the alignments. BAM files were converted to GVCF format using `gatk HaplotypeCaller`. Variant calling was performed using the GVCF files for the analysis cohort. The `gatk GenomicsDBImport` and `gatk GenotypeGVCFs` algorithms were used to generate a master vcf file which was subsequently filtered with `gatk VariantFiltration` (parameters below). The variants that passed filtering were used for PCA analysis.

PCA was performed with `plink258`. To avoid overrepresenting species from the cohort, a core subset of the samples was used for PC 1 and 2 calculations (**Table S2**, column S). Eigenvectors from the whole cohort were then plotted on this PC space and visualized using `plotly`.

Note that GATK  $\geq v4.3.0.0$  and python  $\geq 3.7$  were used for these analyses. Custom python scripts were written to automate and parallelize processing of the samples in the cohorts.

The commands used for variant calling and PCA are the following:

**UBAM Generation:**

```
> gatk FastqToSam --FASTQ <fq1> --FASTQ2 <fq2> --READ_GROUP_NAME <RUNID> --  
TMP_DIR <Temp_dir> --OUTPUT <temp_bam_file> --SAMPLE_NAME <sample_name> --  
LIBRARY_NAME <library_name> --PLATFORM <platform>
```

**Alignment to GT3:**

```
> bwa mem -t <threads> -M -p <GT3_FASTA> <UBAM_file>
```

**Merge metadata from UBAM to BAM:**

```
> gatkMergeBamAlignment -R <GT3_FASTA> --UNMAPPED_BAM <UBAM_file> --  
ALIGNED_BAM <BAM_file>
```

**Generate GVCFS:**

```
> gatk HaplotypeCaller -R <GT3_FASTA> -I <BAM_file> -ERC GVCF -O <GVCF_file>  
-ERC GVCF is used to generate GVCF files that can be used for subsequent joint genotyping
```

**Read Depth Information:**

```
> gatk CountReads -I <BAM_file>
```



528 **Generate GenomicsDB:**

529 > gatk --java-options -Xmx <memory> GenomicsDBImport --genomicsdb-workspace-  
530 path <database\_output\_path> --intervals <interval> --sample-name-map  
531 <cohort\_samples.txt>

532 **Joint Genotyping:**

533 > gatk --java-options -Xmx <memory> GenotypeGVCFs -R <GT3\_FASTA> -V  
534 <path\_to\_database> -O <VCF\_file> --heterozygosity 0.00175 -A  
535 DepthPerAlleleBySample -A Coverage -A GenotypeSummaries -A TandemRepeat -A  
536 StrandBiasBySample -A ReadPosRankSumTest -A AS\_ReadPosRankSumTest -A  
537 AS\_QualByDepth -A AS\_StrandOddsRatio -A AS\_MappingQualityRankSumTest -A  
538 FisherStrand -A QualByDepth RMSMappingQuality -A DepthPerSampleHC -G  
539 StandardAnnotation AS\_StandardAnnotation -G StandardHCAnnotation

540 Note that an average measure of heterozygosity, 0.00175 was used from what was reported in Malinksky *et al.*<sup>29</sup>.

542 **Filtering Variants:**

543 > gatk VariantFiltration -R <GT3\_FASTA> <VCF\_file> -O <Filtered\_VCF> --  
544 filter-name 'allele\_freq' --filter-expression 'AF < 0.05' --filter-name  
545 'inbreeding\_test' --filter-expression 'InbreedingCoeff < -0.6' --filter-name  
546 'depth\_Qual' --filter-expression 'QD < 2.0' --filter-name 'max\_DP' --filter-  
547 expression 'DP > 11000' --filter-name 'min\_DP' --filter-expression 'DP <  
548 7600' --filter-name 'strand\_bias' --filter-expression 'FS > 40.0' --filter-  
549 name 'mapping\_quality' --filter-expression 'MQ < 50.0' --filter-name  
550 'no\_calls' --filter-expression 'NCC > 119' --verbosity ERROR"))

551 Variants were filtered by allele frequency, inbreeding coefficient, quality by depth, depth, fisher's exact  
552 test for strand bias, mapping quality, and by excess missingness according to methods published in  
553 Malinksky *et al.*<sup>29</sup>. We filtered by depth to include variants with a mean depth per sample between 10% and  
554 95% of the distribution of all variant depths.

555 > gatk SelectVariants -V <Filtered\_VCF> --exclude-filtered -O  
556 <pass\_VCF\_file>

557  
558 **PCA:**

559 For PCA, first a subset of the pass\_VCF\_file containing the samples we used to create the principal  
560 components was generated using *bcftools*<sup>59</sup>. Next, the relevant pfiles were generated for both the  
561 pass\_VCF\_file and the subset\_samples\_VCF. Linkage pruning was only performed for the whole genome  
562 and whole chromosome PCAs using the `-indep-pairwise 50 5 0.1` flag and parameters. Since inversions  
563 inherently link together variants, linkage pruning was not used when restricting the analysis to within  
564 inverted regions. A linear scoring system is generated using the subset\_samples\_VCF pfiles. This scoring  
565 system is applied to the genotype matrix for the whole cohort's samples to scale each sample  
566 consistently. The resulting first two eigenvectors per sample from the .sscore file are plotted with *plotly*.

567 > plink2 --vcf <pass\_VCF\_file> --out <whole\_pfiles>  
568 > plink2 --vcf <subset\_samples\_VCF\_file> --out <subset\_pfiles>  
569 > plink2 --pfile <whole\_pfiles> --set-missing-var-ids @:# --make-pgen --out  
570 <whole\_corrected>

571

572 LD Pruning for Whole genome and whole chromosome PCAs:

```
573 > plink2 -pfile <subset_pfiles> --out
574 <subset_samples_ld_pruning_intermediate> --allow-extra-chr --set-missing-
575 var-ids @:#' --indep-pairwise 50 5 0.1
576 > plink2 -pfile <subset_pfiles> --freq counts -pca allele-wts -out
577 <sample_subset_yes_ld_pruning_pca> --allow-extra-chr --set-missing-var-ids
578 @:# --max-alleles 2 --extract <
579 subset_samples_ld_pruning_intermediate.prune.in
580 > plink2 -pfile <whole_corrected> --read-freq <
581 sample_subset_yes_ld_pruning_pca.acount> --score <
582 sample_subset_yes_ld_pruning_pca.eigenval.allele> 2 5 header-read no-mean-
583 imputation variance-standardize --score-col-nums 6-15 -out
584 <projected_pca_yes_ld --allow-extra-chr
```

585

586 No LD Pruning for inverted region PCAs:

```
587 > plink2 --pfile <subset_pfiles> --freq counts -pca allele-wts -out
588 <subset_samples_no_ld_pruning_pca> --allow-extra-chr --set-missing-var-ids
589 @:# --max-alleles 2
590 > plink2 -pfile <whole_corrected> --read-freq <
591 subset_samples_no_ld_pruning_pca.acount> --score'
592 <subset_samples_no_ld_pruning_pca.eigenvec.allele> 2 5 header-read no-mean-
593 imputation variance-standardize --score-col-nums 6-15 --out
594 <projected_PCA_no_ld> --allow-extra-chr
```

595 **Genome alignment.** Genomes were aligned using `minimap2`<sup>60</sup> (v2.22) to align the query genome to the `M_zebra_GT3` reference using default settings. Custom scripts were used to convert the output paf files to dataframe objects, filtering removed alignments that were 1) secondary matches, 2) had less than 30% percent identify, and 3) were less than 10,000 base pairs in length. Further, we excluded any alignments from a contig containing less than 2,000,000 total bp of alignments, to account for repetitive DNA. Alignments were plotted using `seaborn`. We estimated the intervals containing the breakpoints for the inversions using genome alignments of *Aulonocara* sp. 'chitande type north' Nkhata Bay and `M_zebra_GT3` (Table S4).

603 **Phylogenies.** To create phylogenies, `bcftools` was used to filter out non SNV variants from the master vcf file. `Bcftools` was also used to create individual vcf files for each inversion, filtering out SNVs that fell outside of the inverted region. `Vcf2philip`<sup>61</sup> was used to convert the to the phylip format, using the `-m 50` option to filter out SNVs with less than 50 genotyped samples. `Iqtree2`<sup>62</sup> was used to create trees and estimate confidence values for each node using the following options: `-nt', '24', '-mem', '96G', '-v', '--seqtype', 'DNA', '-m', 'GTR+I+G', '-B', '1000'`. Trees were visualized with `iTol`<sup>63</sup> or `ete3`<sup>64</sup>.

610 **Nucleotide diversity analysis.** To analyze the Pixy between individuals of different species, we used `scikit-allel` (v1.3.8)<sup>65</sup> to calculate the genetic difference between each sample using the `allel.pairwise_distance` function and the `citblock` for the metric. We restricted the analysis to benthic animals that carried the inversion and compared these to animals of various ecogroups. We used the `seaborn` function `hist plot` to create density histograms for each of the inversions.

615 **Fst & pedigree.** A vcf file was created for the *Aulonocara* samples listed in **Table S2** (Column I). The PCA  
616 approach (described above) was used to genotype each of the samples for the inversion on 10. For the  $F_{ST}$   
617 analysis, we used `scikit-allel` (v1.3.8) to perform `allel.average_weir_cockerham_fst` with  
618 a window size of 100. For the heterozygosity analysis, we used the  
619 `allel.heterozygosity_observed` function on both the male and female population, subtracting  
620 the average heterozygosity from the male offspring from the female offspring.

## 622 DATA AVAILABILITY

623 All Illumina and PacBio sequencing reads have been deposited to the NCBI Short Read archive at  
624 BioProject PRJNA1112855. The VCF files used in these analyses are available through the Dryad Digital  
625 Repository (TBD).

## 627 REFERENCES

- 628 1. Kirkpatrick, M. How and Why Chromosome Inversions Evolve. *PLOS Biol.* **8**, e1000501 (2010).
- 629 2. Wellenreuther, M. & Bernatchez, L. Eco-Evolutionary Genomics of Chromosomal Inversions. *Trends*  
630 *Ecol. Evol.* **33**, 427–440 (2018).
- 631 3. Berdan, E. L. *et al.* How chromosomal inversions reorient the evolutionary process. *J. Evol. Biol.* **36**,  
632 1761–1782 (2023).
- 633 4. Hoffmann, A. A. & Rieseberg, L. H. Revisiting the Impact of Inversions in Evolution: From Population  
634 Genetic Markers to Drivers of Adaptive Shifts and Speciation? *Annu. Rev. Ecol. Evol. Syst.* **39**, 21–42  
635 (2008).
- 636 5. Andolfatto, P., Depaulis, F. & Navarro, A. Inversion polymorphisms and nucleotide variability in  
637 *Drosophila*. *Genet. Res.* **77**, 1–8 (2001).
- 638 6. Harringmeyer, O. S. & Hoekstra, H. E. Chromosomal inversion polymorphisms shape the genomic  
639 landscape of deer mice. *Nat. Ecol. Evol.* **6**, 1965–1979 (2022).
- 640 7. Huang, K., Andrew, R. L., Owens, G. L., Ostevik, K. L. & Rieseberg, L. H. Multiple chromosomal  
641 inversions contribute to adaptive divergence of a dune sunflower ecotype. *Mol. Ecol.* **29**, 2535–2549  
642 (2020).

- 643 8. Hager, E. R. *et al.* A chromosomal inversion contributes to divergence in multiple traits between deer  
644 mouse ecotypes. *Science* **377**, 399–405 (2022).
- 645 9. Joron, M. *et al.* Chromosomal rearrangements maintain a polymorphic supergene controlling butterfly  
646 mimicry. *Nature* **477**, 203–206 (2011).
- 647 10. Joron, M. *et al.* A Conserved Supergene Locus Controls Colour Pattern Diversity in Heliconius  
648 Butterflies. *PLOS Biol.* **4**, e303 (2006).
- 649 11. Palmer, D. H. & Kronforst, M. R. A shared genetic basis of mimicry across swallowtail butterflies  
650 points to ancestral co-option of doublesex. *Nat. Commun.* **11**, 6 (2020).
- 651 12. Jones, F. C. *et al.* The genomic basis of adaptive evolution in threespine sticklebacks. *Nature* **484**,  
652 55–61 (2012).
- 653 13. Lamichhaney, S. *et al.* Structural genomic changes underlie alternative reproductive strategies in  
654 the ruff (*Philomachus pugnax*). *Nat. Genet.* **48**, 84–88 (2016).
- 655 14. Küpper, C. *et al.* A supergene determines highly divergent male reproductive morphs in the ruff.  
656 *Nat. Genet.* **48**, 79–83 (2016).
- 657 15. Wang, J. *et al.* A Y-like social chromosome causes alternative colony organization in fire ants.  
658 *Nature* **493**, 664–668 (2013).
- 659 16. Tuttle, E. M. *et al.* Divergence and Functional Degradation of a Sex Chromosome-like Supergene.  
660 *Curr. Biol. CB* **26**, 344–350 (2016).
- 661 17. Hughes, J. F. *et al.* Chimpanzee and human Y chromosomes are remarkably divergent in structure  
662 and gene content. *Nature* **463**, 536–539 (2010).
- 663 18. Fishman, L., Stathos, A., Beardsley, P. M., Williams, C. F. & Hill, J. P. CHROMOSOMAL  
664 REARRANGEMENTS AND THE GENETICS OF REPRODUCTIVE BARRIERS IN MIMULUS (MONKEY  
665 FLOWERS). *Evolution* **67**, 2547–2560 (2013).

- 666 19. Noor, M. A., Grams, K. L., Bertucci, L. A. & Reiland, J. Chromosomal inversions and the  
667 reproductive isolation of species. *Proc. Natl. Acad. Sci. U. S. A.* **98**, 12084–12088 (2001).
- 668 20. Fuller, Z. L., Koury, S. A., Phadnis, N. & Schaeffer, S. W. How chromosomal rearrangements shape  
669 adaptation and speciation: Case studies in *Drosophila pseudoobscura* and its sibling species  
670 *Drosophila persimilis*. *Mol. Ecol.* **28**, 1283–1301 (2019).
- 671 21. Fuller, Z. L., Leonard, C. J., Young, R. E., Schaeffer, S. W. & Phadnis, N. Ancestral polymorphisms  
672 explain the role of chromosomal inversions in speciation. *PLOS Genet.* **14**, e1007526 (2018).
- 673 22. Kirkpatrick, M. & Barton, N. Chromosome Inversions, Local Adaptation and Speciation. *Genetics*  
674 **173**, 419–434 (2006).
- 675 23. Lowry, D. B. & Willis, J. H. A Widespread Chromosomal Inversion Polymorphism Contributes to a  
676 Major Life-History Transition, Local Adaptation, and Reproductive Isolation. *PLoS Biol.* **8**, e1000500  
677 (2010).
- 678 24. Trickett, A. J. & Butlin, R. K. Recombination suppressors and the evolution of new species.  
679 *Heredity* **73**, 339–345 (1994).
- 680 25. Kocher, T. D. Adaptive evolution and explosive speciation: the cichlid fish model. *Nat. Rev. Genet.*  
681 **5**, 288–298 (2004).
- 682 26. Santos, M. E., Lopes, J. F. & Kratochwil, C. F. East African cichlid fishes. *EvoDevo* **14**, 1 (2023).
- 683 27. Svardal, H., Salzburger, W. & Malinsky, M. Genetic Variation and Hybridization in Evolutionary  
684 Radiations of Cichlid Fishes. *Annu. Rev. Anim. Biosci.* **9**, 55–79 (2021).
- 685 28. Johnson, Z. V. *et al.* Cellular profiling of a recently-evolved social behavior in cichlid fishes. *Nat.*  
686 *Commun.* **14**, 4891 (2023).
- 687 29. Malinsky, M. *et al.* Whole-genome sequences of Malawi cichlids reveal multiple radiations  
688 interconnected by gene flow. *Nat. Ecol. Evol.* **2**, 1940–1955 (2018).

- 689 30. Patil, C. *et al.* Genome-enabled discovery of evolutionary divergence in brains and behavior. *Sci.*  
690 *Rep.* **11**, 13016 (2021).
- 691 31. Todd Strelman, J. & Danley, P. D. The stages of vertebrate evolutionary radiation. *Trends Ecol.*  
692 *Evol.* **18**, 126–131 (2003).
- 693 32. Maan, M. E. & Sefc, K. M. Colour variation in cichlid fish: Developmental mechanisms, selective  
694 pressures and evolutionary consequences. *Semin. Cell Dev. Biol.* **24**, 516–528 (2013).
- 695 33. Konings, A. Fishes, as well as birds, build bowers.
- 696 34. Martin, C. & Genner, M. A role for male bower size as an intrasexual signal in a Lake Malawi cichlid  
697 fish. (2009) doi:10.1163/156853908X396836.
- 698 35. Loh, Y.-H. E. *et al.* Origins of Shared Genetic Variation in African Cichlids. *Mol. Biol. Evol.* **30**, 906–  
699 917 (2013).
- 700 36. Meier, J. I. *et al.* Cycles of fusion and fission enabled rapid parallel adaptive radiations in African  
701 cichlids. *Science* **381**, eade2833 (2023).
- 702 37. Keller, I. *et al.* Population genomic signatures of divergent adaptation, gene flow and hybrid  
703 speciation in the rapid radiation of Lake Victoria cichlid fishes. *Mol. Ecol.* **22**, 2848–2863 (2013).
- 704 38. Brawand, D. *et al.* The genomic substrate for adaptive radiation in African cichlid fish. *Nature* **513**,  
705 375–381 (2014).
- 706 39. Conte, M. A. *et al.* Chromosome-scale assemblies reveal the structural evolution of African cichlid  
707 genomes. *GigaScience* **8**, giz030 (2019).
- 708 40. Lam, E. T. *et al.* Genome mapping on nanochannel arrays for structural variation analysis and  
709 sequence assembly. *Nat. Biotechnol.* **30**, 10.1038/nbt.2303 (2012).
- 710 41. Ciezarek, A. G. *et al.* Ancient and Recent Hybridization in the Oreochromis Cichlid Fishes. *Mol.*  
711 *Biol. Evol.* **41**, msae116 (2024).

- 712 42. Nowling, R. J., Manke, K. R. & Emrich, S. J. Detecting inversions with PCA in the presence of  
713 population structure. *PLoS ONE* **15**, e0240429 (2020).
- 714 43. York, R. A. *et al.* Behavior-dependent cis regulation reveals genes and pathways associated with  
715 bower building in cichlid fishes. *Proc. Natl. Acad. Sci.* **115**, E11081–E11090 (2018).
- 716 44. Kratochwil, C. F., Liang, Y., Urban, S., Torres-Dowdall, J. & Meyer, A. Evolutionary Dynamics of  
717 Structural Variation at a Key Locus for Color Pattern Diversification in Cichlid Fishes. *Genome Biol.*  
718 *Evol.* **11**, 3452–3465 (2019).
- 719 45. Feulner, P. G. D., Schwarzer, J., Haesler, M. P., Meier, J. I. & Seehausen, O. A Dense Linkage Map of  
720 Lake Victoria Cichlids Improved the Pundamilia Genome Assembly and Revealed a Major QTL for Sex-  
721 Determination. *G3 GenesGenomesGenetics* **8**, 2411–2420 (2018).
- 722 46. Johnson, N. A. & Lachance, J. The genetics of sex chromosomes: evolution and implications for  
723 hybrid incompatibility. *Ann. N. Y. Acad. Sci.* **1256**, E1-22 (2012).
- 724 47. Natri, H. M., Merilä, J. & Shikano, T. The evolution of sex determination associated with a  
725 chromosomal inversion. *Nat. Commun.* **10**, 145 (2019).
- 726 48. Kondrashov, A. S. & Mina, M. V. Sympatric speciation: when is it possible? *Biol. J. Linn. Soc.* **27**,  
727 201–223 (1986).
- 728 49. Danley, P. D. & Kocher, T. D. Speciation in rapidly diverging systems: lessons from Lake Malawi.  
729 *Mol. Ecol.* **10**, 1075–1086 (2001).
- 730 50. Parnell, N. F. & Streelman, J. T. Genetic interactions controlling sex and color establish the  
731 potential for sexual conflict in Lake Malawi cichlid fishes. *Heredity* **110**, 239–246 (2013).
- 732 51. Behrens, K. A., Koblmüller, S. & Kocher, T. D. Diversity of Sex Chromosomes in Vertebrates: Six  
733 Novel Sex Chromosomes in Basal Haplochromines (Teleostei: Cichlidae). *Genome Biol. Evol.* **16**,  
734 evae152 (2024).

- 735 52. Blumer, L. M. *et al.* Introgression dynamics of sex-linked chromosomal inversions shape the  
736 Malawi cichlid adaptive radiation. Preprint at <https://doi.org/10.1101/2024.07.28.605452> (2024).
- 737 53. Schelkunov, M. I. Mabs, a suite of tools for gene-informed genome assembly. *BMC Bioinformatics*  
738 **24**, 377 (2023).
- 739 54. Chen, Y., Zhang, Y., Wang, A. Y., Gao, M. & Chong, Z. Accurate long-read de novo assembly  
740 evaluation with Inspector. *Genome Biol.* **22**, 312 (2021).
- 741 55. Cabanettes, F. & Klopp, C. D-GENIES: dot plot large genomes in an interactive, efficient and  
742 simple way. *PeerJ* **6**, e4958 (2018).
- 743 56. Poplin, R. *et al.* Scaling accurate genetic variant discovery to tens of thousands of samples.  
744 201178 Preprint at <https://doi.org/10.1101/201178> (2018).
- 745 57. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform.  
746 *Bioinformatics* **25**, 1754–1760 (2009).
- 747 58. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets.  
748 *GigaScience* **4**, s13742-015-0047–8 (2015).
- 749 59. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and  
750 population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987–2993  
751 (2011).
- 752 60. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100  
753 (2018).
- 754 61. Ortiz, E. M. vcf2phylip v2.0: convert a VCF matrix into several matrix formats for phylogenetic  
755 analysis. Zenodo <https://doi.org/10.5281/zenodo.2540861> (2019).
- 756 62. Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the  
757 Genomic Era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).



- 758 63. Letunic, I. & Bork, P. Interactive Tree of Life (iTOL) v6: recent updates to the phylogenetic tree  
759 display and annotation tool. *Nucleic Acids Res.* **52**, W78–W82 (2024).
- 760 64. Huerta-Cepas, J., Serra, F. & Bork, P. ETE 3: Reconstruction, Analysis, and Visualization of  
761 Phylogenomic Data. *Mol. Biol. Evol.* **33**, 1635–1638 (2016).
- 762 65. Miles, A. *et al.* cggh/scikit-allele: v1.3.13. Zenodo <https://doi.org/10.5281/zenodo.13772087> (2024).
- 763