

Whole genome sequences of Malawi cichlids reveal multiple radiations interconnected by gene flow

Malinsky, Milan; Svardal, Hannes; Tyers, Alexandra; Miska, Eric.A.; Genner, Martin J.; Turner, George; Durbin, Richard

Nature Ecology and Evolution

DOI:
[10.1038/s41559-018-0717-x](https://doi.org/10.1038/s41559-018-0717-x)

Published: 01/12/2018

Peer reviewed version

[Cyswllt i'r cyhoeddiad / Link to publication](#)

Dyfyniad o'r fersiwn a gyhoeddwyd / Citation for published version (APA):
Malinsky, M., Svardal, H., Tyers, A., Miska, E. A., Genner, M. J., Turner, G., & Durbin, R. (2018). Whole genome sequences of Malawi cichlids reveal multiple radiations interconnected by gene flow. *Nature Ecology and Evolution*, 2, 1940-1955. <https://doi.org/10.1038/s41559-018-0717-x>

Hawliau Cyffredinol / General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

1 **Title: Whole genome sequences of Malawi cichlids reveal multiple radiations**
2 **interconnected by gene flow**

3 **Authors:** Milan Malinsky^{1,2,†*}, Hannes Svardal^{1,3,4,5,†}, Alexandra M. Tyers⁶, Eric A.
4 Miska^{1,3,7}, Martin J. Genner⁸, George F. Turner⁶, and Richard Durbin^{1,3*}

5 **Affiliations:**

6 ¹Wellcome Sanger Institute, Cambridge, CB10 1SA, UK.

7 ²Zoological Institute, University of Basel, 4051 Basel, Switzerland.

8 ³Department of Genetics, University of Cambridge, Cambridge CB2 3EH, UK.

9 ⁴Department of Biology, University of Antwerp, 2020 Antwerp, Belgium.

10 ⁵Naturalis Biodiversity Center, 2300 RA Leiden, The Netherlands.

11 ⁶School of Natural Sciences, Bangor University, Bangor, LL57 2UW, UK.

12 ⁷Gurdon Institute, University of Cambridge, Cambridge, CB2 1QN, UK.

13 ⁸School of Biological Sciences, University of Bristol, Bristol BS8 1TQ, UK.

14 *Correspondence to: millanek@gmail.com (MM), rd@sanger.ac.uk (RD).

15 †These authors contributed equally to this work.

16
17 **Abstract:** The hundreds of cichlid fish species in Lake Malawi constitute the most extensive
18 recent vertebrate adaptive radiation. Here we characterize its genomic diversity by sequencing
19 134 individuals covering 73 species across all major lineages. Average sequence divergence
20 between species pairs is only 0.1-0.25%. These divergence values overlap diversity within
21 species, with 82% of heterozygosity shared between species. Phylogenetic analyses suggest that
22 diversification initially proceeded by serial branching from a generalist *Astatotilapia*-like
23 ancestor. However, no single species tree adequately represents all species relationships, with
24 evidence for substantial gene flow at multiple times. Common signatures of selection on visual
25 and oxygen transport genes shared by distantly related deep water species point to both adaptive
26 introgression and independent selection. These findings enhance our understanding of genomic
27 processes underlying rapid species diversification, and provide a platform for future genetic
28 analysis of the Malawi radiation.

30 **Main Text:** The formation of every lake or island represents a fresh opportunity for
31 colonization, proliferation and diversification of living forms. In some cases, the ecological
32 opportunities presented by underutilized habitats facilitate adaptive radiation - rapid and
33 extensive diversification of the descendants of the colonizing lineages¹⁻³. Adaptive radiations are
34 thus exquisite examples of the power of natural selection, as seen for example in Darwin's
35 finches in the Galapagos^{4,5}, Anolis lizards of the Caribbean⁶ and in East African cichlid fishes^{7,8}.

36 Cichlids are one of the most species-rich and diverse families of vertebrates, and nowhere are
37 their radiations more spectacular than in the Great Lakes of East Africa: Malawi, Tanganyika,
38 and Victoria², each of which contains several hundred endemic species, with the largest number
39 in Lake Malawi⁹. Molecular genetic studies have made major contributions to reconstructing the
40 evolutionary histories of these adaptive radiations, especially in terms of the relationships
41 between the lakes^{10,11}, between some major lineages in Lake Tanganyika¹², and in describing the
42 role of hybridization in the origins of the Lake Victoria radiation¹³. However, the task of
43 reconstructing within-lake relationships remains challenging due both to retention of large
44 amounts of ancestral genetic polymorphism (i.e. incomplete lineage sorting) and gene flow
45 between taxa^{12,14-18}.

46 Initial genome assemblies of cichlids from East Africa suggest that an increased rate of gene
47 duplication together with accelerated evolution of some regulatory elements and protein coding
48 genes may have contributed to the radiations¹¹. However, understanding of the genomic
49 mechanisms contributing to adaptive radiations is still in its infancy³.

50 Here we provide an overview of and insights into the genomic signatures of the haplochromine
51 cichlid radiation of Lake Malawi. The species that comprise the radiation can be divided into
52 seven groups with differing ecology and morphology (see Supplementary Note): 1) the rock-
53 dwelling 'mbuna'; 2) *Rhamphochromis* - typically midwater pelagic piscivores; 3) *Diplotaxodon*
54 - typically deep-water pelagic zooplanktivores and piscivores; 4) deep-water and twilight feeding
55 benthic species; 5) 'utaka' feeding on zooplankton in the water column but breeding on or near
56 the lake bottom (here utaka corresponds to the genus *Copadichromis*); 6) a diverse group of
57 benthic species, mainly found in shallow non-rocky habitats. In addition, *Astatotilapia calliptera*
58 is a closely related generalist that inhabits shallow weedy margins of Lake Malawi, and other

59 lakes and rivers in the catchment, as well as river systems to the east and south of the Lake
60 Malawi catchment. This division into seven groups has been partially supported by previous
61 molecular phylogenies based on mtDNA and amplified fragment length polymorphism (AFLP)
62 data¹⁸⁻²⁰. However, published phylogenies show numerous inconsistencies and, in particular, the
63 question of whether the groups are genetically separate remained unanswered.

64 To characterize the genetic diversity, species relationships, and signatures of selection across the
65 whole radiation, we obtained Illumina whole-genome sequence data from 134 individuals of 73
66 species distributed broadly across the seven groups (Fig. 1a; Supplementary Note). This includes
67 102 individuals at ~15× coverage and 32 additional individuals at ~6× (Supplementary Table 1).

68 **Results**

69 **Low genetic diversity and species divergence**

70 Sequence data were aligned to and variants called against a *Metriaclima zebra* reference
71 genome¹¹. Average divergence from the reference was 0.19% to 0.27% (Supplementary Fig. 1).
72 After filtering and variant refinement, we obtained 30.6 million variants of which 27.1 million
73 were single nucleotide polymorphisms (SNPs) and the rest were short insertions and deletions.
74 All the following analyses are based on biallelic SNPs.

75 To estimate nucleotide diversity (π) within the species, we measured the frequency of
76 heterozygous sites in each individual. The estimates are distributed within a relatively narrow
77 range between 0.7 and 1.8×10^{-3} per bp (Fig. 1b). The mean π estimate of 1.2×10^{-3} per bp is at the
78 low end of values found in other animals²¹. There does not appear to be a relationship between π
79 and the rate of speciation: individuals in the species-rich mbuna and shallow benthic groups
80 show levels of π comparable to the relatively species-poor utaka, *Diplotaxodon*, and
81 *Rhamphochromis* (Supplementary Fig. 1).

82 Despite their extensive phenotypic differentiation, species within the Lake Malawi radiation are
83 genetically closely related^{22,23}. However, genome-wide genetic divergence has never been
84 quantified. We calculated the average pairwise sequence differences (d_{XY}) between species and
85 compared d_{XY} against heterozygosity, finding that the two distributions partially overlap (Fig.
86 1b). Thus, the sequence divergence within a single diploid individual is sometimes higher than
87 the divergence between two distinct species. The average d_{XY} is 2.0×10^{-3} with a range between

88 1.0 and 2.4×10^{-3} per bp. The maximum d_{XY} is therefore approximately one fifth of the
89 divergence between human and chimpanzee²⁴. In addition to the low ratio of divergence to
90 diversity, most genetic variation is shared between species. On average both alleles are observed
91 in other species for 82% of heterozygous sites within individuals, consistent with the expected
92 and previously observed high levels of incomplete lineage sorting (ILS)²³. Supplementary Fig. 2
93 shows d_{XY} and F_{ST} values for comparisons between the seven eco-morphological groups and
94 Supplementary Fig. 3 shows patterns of linkage disequilibrium across the radiation, within
95 groups, and within individual species.

96 **Low per-generation mutation rate**

97 It has been suggested that the species richness and morphological diversity of teleosts in general
98 and of cichlids in particular might be explained by elevated mutation rates compared to other
99 vertebrates^{25,26}. To obtain a direct estimate of the per-generation mutation rate, we reared
100 offspring of three species from three different Lake Malawi groups (*A. calliptera*, *Aulonocara*
101 *stuartgranti* and *Lethrinops lethrinus*). We sequenced both parents and one offspring of each to
102 high coverage (40x), applied stringent quality filtering, and counted variants present in each
103 offspring but absent in both its parents (Supplementary Fig. 4). There was no evidence for
104 significant difference in mutation rates between species. The overall mutation rate (μ) was
105 estimated at 3.5×10^{-9} (95%CI: 1.6×10^{-9} to 4.6×10^{-9}) per bp per generation, approximately three
106 to four times lower than in human²⁷, although, given much shorter mean generation times, the
107 per-year rate is still expected to be higher in cichlids than in humans. We note that Recknagel et
108 al.²⁶ obtained a much higher mutation rate estimate (6.6×10^{-8} per bp per generation) in Midas
109 cichlids, but from relatively low depth RADseq data that may have made accurate verification
110 more difficult. We also note that our per generation rate estimate, although low, is still higher
111 than the lowest μ estimate in vertebrates: 2×10^{-9} per bp per generation recently reported for
112 Atlantic herring²⁸. By combining our mutation rate with nucleotide diversity (π) values, we
113 estimate the long term effective population sizes (N_e) to be in the range of approximately 50,000
114 to 130,000 breeding individuals (with $N_e = \pi/4\mu$).

115
116
117
118
119

120 **Genome data support for eco-morphological groupings**

121 Principal component analysis (PCA) of the whole-genome genotype data generally separates the
122 major eco-morphological groups (Fig. 1c). The most notable exceptions to this are (1) the utaka,
123 for which some species cluster more closely with deep benthics and others with shallow
124 benthics, and (2) two species of the genus *Aulonocara*, *A. stuartgranti* and *A. steveni*, which are
125 located between the shallow and deep benthic groups. Although these have enlarged lateral line
126 sensory apparatus like many deep benthic species including other *Aulonocara*, they are typically
127 found in shallower water²⁹. Another interesting pattern in the PCA plot is that the utaka and
128 benthic samples are often spread along principal component (PC) axes (Fig. 1c, Supplementary
129 Fig. 5), a pattern typical for admixed populations (e.g. ref. 30). Along the two main PCs, the
130 deeper water benthic species extend towards the deep water *Diplotaxodon*, an observation we
131 will return to in the context of gene flow and shared mechanisms of depth adaptation.

132 To further verify the consistency of group assignments, we tested whether pairs of species from
133 the same group always share more derived alleles with each other than with any species from
134 other groups. Group assignments were again supported, except for the four species also
135 highlighted in the PCA: the two shallow-living *Aulonocara* are closer to shallow benthics than to
136 deep benthics in 71% and 82% of tests respectively when comparing these alternatives, and
137 *Copadichromis trimaculatus* is closer to shallow benthics than to utaka in 58% of the
138 comparisons. *Copadichromis cf. trewavasae* always clustered with shallow benthics; therefore,
139 we treat it as a member of the shallow benthic group throughout the remainder of this
140 manuscript. With the three intermediate samples removed and *C. cf. trewavasae* reassigned, all
141 other species showed 100% consistency with their group assignment.

142 **Allele sharing inconsistent with tree-like relationships**

143 The above observations suggest that some species may be genetic intermediates between well-
144 defined groups, consistent with previous studies which have suggested that hybridization and
145 introgression subsequent to initial separation of species may have played a significant role in
146 cichlid radiations, including in Lakes Tanganyika^{12,14-16} and Malawi^{18,20}. Where this happens,
147 there is no single tree relating the species.

148 To assess the overall extent of violation of tree-like species relationships, we calculated
149 Patterson's D statistic (ABBA-BABA test)^{31,32} for all possible trios of Lake Malawi species,
150 without assuming any *a priori* knowledge of their relationships. *N. brichardi* from Lake
151 Tanganyika was always used as the outgroup. The test statistic D_{\min} is the minimum absolute
152 value of Patterson's D for each trio, across all possible tree topologies. Therefore, a significantly
153 positive D_{\min} score signifies that the sharing of derived alleles between the three species is
154 inconsistent with a single species tree relating them, even in the presence of incomplete lineage
155 sorting.

156 Overall, 62% of trios (75,616 out of 121,485) have significantly positive D_{\min} score (Holm-
157 Bonferroni FWER < 0.01). The D_{\min} values are not independent: for example, a single gene flow
158 event between ancestral lineages can affect multiple contemporary species and thus more trios
159 than a more recent gene flow event would. However, tree violations are numerous and pervasive
160 throughout the dataset, within all the major groups and also between groups (Fig. 2a), revealing
161 reticulate evolution at multiple levels. Therefore, phylogenetic trees alone cannot fully describe
162 the evolutionary relationships of Lake Malawi cichlids.

163

164 **Phylogenetic framework**

165 Despite no tree giving a complete and accurate picture of the relationships between species,
166 standard phylogenetic approaches are useful to provide a framework for discussion. To obtain an
167 initial picture we divided the genome into 2543 non-overlapping windows, each comprising
168 8000 SNPs (average size: 274kb) and constructed a Maximum Likelihood (ML) phylogeny
169 separately for the full sequences within each window, obtaining trees with 2542 different
170 topologies. We also calculated the maximum clade credibility (MCC) summary tree³³ and an ML
171 phylogeny based on the full mtDNA genome (Fig. 1c and Supplementary Fig. 6).

172 We next applied a range of further phylogenomic methods which are known to be robust to
173 incomplete lineage sorting. These included three multispecies coalescent methods^{34,35}: the
174 Bayesian SNAPP³⁶ (with a subset of 48,922 unlinked SNPs in 12 individuals representing the
175 eco-morphological groups), the algebraic method SVDquartets^{37,38}, which allows for site-specific
176 rate variation and is robust to gene-flow between sister taxa³⁹, and the summary method

177 ASTRAL^{40,41}, using the 2543 local ML trees that were described above as input. We also built a
178 whole genome Neighbour-Joining (NJ) tree using the Dasarathy et al. algorithm, which has been
179 shown to be a statistically consistent and accurate species tree estimator under ILS^{42,43}. The
180 above methods have also been applied to datasets where the individuals that are genetically
181 intermediate between eco-morphological groups (*C. trimaculatus*, *A. stuartgranti*, and *A. steveni*)
182 have been removed, thus likely reducing the extent of violation of the multispecies coalescent
183 model.

184 Despite extensive variation among the 2543 individual ML trees (at least in part attributable to
185 ILS), and, to a lesser extent, variation between the different genome-wide phylogenetic methods,
186 there is some general consensus (Fig. 2c and Supplementary Figs. 6-10). Except for the three
187 previously identified intermediate species, individuals from within each of the previously
188 identified eco-morphological groups cluster together in all the whole genome phylogenies,
189 forming well supported reciprocally monophyletic groups. The pelagic *Diploaxodon* and
190 *Rhamphochromis* together form a sister group to the rest of the radiation, except in the all-sample
191 MCC and SVDquartets phylogenies. Perhaps surprisingly, all the methods place the generalist *A.*
192 *calliptera* as the sister taxon to the specialized rocky-shore mbuna group in a position that is
193 nested within the Lake Malawi radiation. On a finer scale, many similarities between the
194 resulting phylogenies reflect features of previous taxonomic assignment, but some currently-
195 recognized genera are always polyphyletic, including *Placidochromis*, *Lethrinops*, and
196 *Mylochromis*.

197 The mtDNA phylogeny is an outlier, substantially different from all the whole-genome
198 phylogenies and also from the majority of the local ML trees (Fig. 2b,c and Supplementary Figs.
199 6 and 11). Discordances between mtDNA and nuclear phylogenies in Lake Malawi have been
200 reported previously and interpreted as a signature of past hybridization events^{18,20}. However, as
201 we discuss below, some of these previously suggested hybridization events are not reflected in
202 the whole genome data. Indeed, large discrepancies between mitochondrial and nuclear
203 phylogenies have been shown in many other systems, reflecting both that mtDNA as a single
204 locus is not expected to reflect the consensus under ILS, and high incidence of mitochondrial
205 selection⁴⁴⁻⁴⁶. This underlines the importance of evaluating species relationships in the Lake
206 Malawi radiation from a genome-wide perspective.

207

208 **Specific signals of introgression**

209 We applied a variety of methods to identify the species and groups whose relationships violate
210 the framework trees described in the previous section. First, we contrasted the pairwise genetic
211 distances used to produce the NJ tree against the distances between samples along the tree
212 branches, calculating the residuals (Supplementary Fig. 12). If the tree captured all the genetic
213 relationships in our sample perfectly, the residuals would all be zero. However, as expected in
214 the light of the D_{\min} analysis above, we found numerous differences, affecting both groups of
215 species and individual species, with some standout cases. Among the strongest signals on
216 individual species, in addition to the previously discussed *C. trimaculatus*, we can see that 1)
217 *Placidochromis cf. longimanus* is genetically closer to the deep benthic clade and to a subset of
218 the shallow benthic (mainly *Lethrinops* species) than the tree suggests; and 2) our sample of
219 *Otopharynx tetrastigma* (from Lake Ilamba) is much closer to *A. calliptera* (especially to the
220 sample from Lake Kingiri, only 3.2km away) than expected from the tree.

221 Second, the sharing of long haplotypes between otherwise distantly related species is an
222 indication of recent admixture or introgression. To investigate this type of gene flow signature,
223 we used the chromopainter software package⁴⁷ and calculated the ‘coancestry matrix’ of all
224 species - a summary of nearest neighbour (therefore recent) haplotype relationships. The Lake
225 Ilamba *O. tetrastigma* and Lake Kingiri *A. calliptera* also stand out in this analysis, showing a
226 strong signature of recent gene flow between individual species from distinct eco-morphological
227 groups (Supplementary Fig. 13). The other tree-violation signatures described above are also
228 visible on the haplotype sharing level but are less pronounced, consistent with being older events
229 involving the common ancestors of multiple present-day species. However, the chromopainter
230 results indicate additional recent introgression events (e.g. the utaka *C. virginalis* with
231 *Diplotaxodon*; more highlighted in Supplementary Fig. 13). Furthermore, the clustering based on
232 recent co-ancestry is different from all phylogenetic trees: in particular a number of shallow
233 benthics, including *P. cf. longimanus*, cluster next to the deep benthics.

234 Third, we used the f_4 admixture ratio^{31,32,48} (f statistic; closely related to Patterson's D),
235 computing $f(A,B;C,O)$ for all groups of species that fit the relationships ((A, B), C) in the
236 ASTRAL* tree (Supplementary Fig. 7), with the outgroup fixed as *N. brichardi*. When elevated

237 due to introgression, the f statistic is expected to be linear in relation to the proportion of
238 introgressed material. The ASTRAL* tree has the lowest mean topological distance to all the
239 other trees, and excludes the three species with intermediate group assignment, a choice made
240 here because we were interested in identifying additional signals beyond the admixed status of *A.*
241 *stuartgranti*, *A. steveni*, and *C. trimaculatus*. Out of the 164,320 computed f statistics, 97,889
242 were significant at FWER < 0.001.

243 As in the case of D_{\min} , a single gene flow event can lead to multiple significant f statistics.
244 Noting that the values for different combinations of ((A, B), C) groups are not independent as
245 soon as they share branches on the tree, we sought to obtain branch-specific estimates of excess
246 allele sharing that would be less correlated. Building on the logic employed to understand
247 correlated gene flow signals in ref. 49, we developed “ f -branch” or $f_b(C)$: a summary of f scores
248 that, on a given tree, captures excess allele sharing between a species C and a branch b compared
249 to the sister branch of b (Methods). Therefore, an $f_b(C)$ score is specific to the branch b (on the y-
250 axis in Fig. 3), but a single introgression event can still lead to significant $f_b(C)$ ’s across multiple
251 related C . There were 11,158 $f_b(C)$ scores of which 1,421 were significantly elevated at
252 FWER<0.001 (Supplementary Fig. 14), and 238 scores were larger than 3% (the value inferred
253 for human-Neanderthal introgression in ref. 31). The majority of nodes in the tree are affected:
254 92 of the 158 branches in the phylogeny show significant excess allele sharing with at least one
255 other species C (Fig. 3).

256 Overall, the highest $f_b(C)$ (14.2%) is between the ancestor of the two sampled *Ctenopharynx*
257 species from the shallow benthic group and the utaka *Copadichromis virginalis* (Fig. 3). Notably,
258 *Ctenopharynx* species, particularly *C. intermedius* and *C. pictus* have very large numbers of long
259 slender gill rakers, a feature shared with *Copadichromis* species, and believed to be related to a
260 diet of small invertebrates⁵⁰. Several other benthic lineages also share excess alleles with *C.*
261 *virginalis* at a lower level. Next, the significantly elevated $f_b(C)$ scores between the shallow and
262 the deep benthic lineages suggest that genetic exchanges between these two groups go beyond
263 the clearly admixed shallow-living *Aulonacara* (not included in this analysis). The f -branch
264 signals between *O. tetrastigma* and *A. calliptera* Kingiri are observed in both directions – *A.*
265 *calliptera* Kingiri with shallow benthics (and most strongly *O. tetrastigma*) and *O. tetrastigma*
266 with *A. calliptera* (most strongly *A. calliptera* Kingiri), suggesting bi-directional introgression.

267 At the level of the major eco-morphological groups, the strongest signal indicates that the
268 ancestral lineage of benthics and utaka shares excess derived alleles with *Diplotaxodon* and, to a
269 lesser degree, *Rhamphochromis*, as previously suggested by the PCA plot (Fig 1c). Furthermore,
270 there is evidence for additional ancestry from the pelagic groups in utaka, which could be
271 explained either by an additional, more recent, gene flow event or by differential fixation of
272 introgressed material, possibly due to selection. Reciprocally, *Diplotaxodon* shares excess
273 derived alleles (relative to *Rhamphochromis*) with utaka and deep benthics, as does
274 *Rhamphochromis* with mbuna and *A. calliptera*. Furthermore, mbuna show excess allele sharing
275 (relative to *A. calliptera*) with *Diplotaxodon* and *Rhamphochromis*. (Fig. 3) On the other hand,
276 while ref. 18 suggested gene flow between the deep benthic and mbuna groups on the basis of a
277 discrepancy between mtDNA and nuclear phylogenies, our genome-wide analysis did not find
278 any signal of substantial genetic exchange between these groups.

279 The f statistic tests are robust to the occurrence of incomplete lineage sorting, in the sense that
280 ILS alone cannot generate a significant test result³². We note, however, that pronounced
281 population structure within ancestral species, coupled with rapid succession of speciation events,
282 can also substantially violate the assumptions of a strictly bifurcating species tree and lead to
283 significantly elevated f scores^{32,51}. This needs to be taken into account when interpreting non-
284 treelike relationships, for example among major groups early in the radiation. However, in cases
285 of excess allele sharing between ‘distant’ lineages that are separated by multiple speciation
286 events, ancestral population structure would have needed to segregate through these speciation
287 events without affecting sister lineages, a scenario that is not credible in general. Therefore, we
288 suggest that there is strong evidence for multiple cross-species gene flow events. Additionally,
289 simulations suggest that, compared with treemix⁵², $f_b(C)$ is robust to misspecification of the
290 initial tree (Supplementary Note).

291 Overall, the NJ tree residuals, the haplotype sharing patterns, and the many elevated $f_b(C)$ scores
292 paint a consistent picture. They confirm the extensive violations of the bifurcating species tree
293 model initially revealed by the D_{\min} analysis, and suggest many independent gene flow events at
294 different times during the evolutionary history of the adaptive radiation.

295

296 **Origins of the radiation**

297 The generalist *Astatotilapia calliptera* has been referred to as the 'prototype' for the endemic
298 Lake Malawi cichlids^{29,53}, and discussions concerning the origin of the radiation often centre on
299 ascertaining its relationship to the Malawi species^{20,54}. Previous phylogenetic analyses, using
300 mtDNA and small numbers of nuclear markers, showed inconsistencies in this respect^{18,20,54}. In
301 contrast, our whole genome data indicated a clear and consistent position of the Lake Malawi
302 catchment *A. calliptera* as a sister group to the mbuna, in agreement with the nuclear DNA
303 phylogeny in ref. 18. While it is not certain whether the 320 remaining mbuna species form a
304 monophyletic group with the eight species we used here, the eight species represent the majority
305 of the genera of mbuna and therefore are likely to be representative of much of the genetic
306 diversity within the group.

307 To explore the origins of the Lake Malawi radiation in greater detail, we obtained 24 additional
308 *Astatotilapia* whole genome sequences from outside of Lake Malawi: five *A. calliptera* from
309 Indian Ocean catchments (IOC), thus covering most of its geographical distribution, and 19
310 individuals from seven other *Astatotilapia* species (Supplementary Table 2). We generated new
311 variant calls (Supplementary Methods) and first constructed a NJ tree, finding that all the *A.*
312 *calliptera* (including IOC) cluster as a single group nested at the same place within the radiation,
313 whereas the other *Astatotilapia* species branched off well before the lake radiation (Fig. 4a,b,c).
314 All *A. calliptera* individuals cluster by geography (Fig. 4b,c), except for the specimen from
315 crater Lake Kingiri, whose position in the tree is likely a result of the admixture signals with *O.*
316 *tetrastigma*. Indeed, a NJ tree built only with *A. calliptera* samples (Supplementary Fig. 17)
317 places the Kingiri individual according to geography with the specimens from the nearby crater
318 Lake Massoko and Mbaka River.

319 Applying the same logic as above, we tested whether the position of the *A. calliptera* group in
320 the NJ tree changes when the tree is built without mbuna (as would be expected if *A. calliptera*
321 were affected by hybridization with mbuna). We found that the position of *A. calliptera* is
322 unaffected (Supplementary Fig. 18), suggesting that the nested position is not due to later
323 hybridization. The *f* statistics in Fig. 3 further support this, because the signals involving the
324 whole mbuna or *A. calliptera* groups are modest and do not suggest erroneous placement of these
325 groups in all phylogenetic analyses. Furthermore, the nested position of *A. calliptera* is also

326 supported by the vast majority of the genome. Searching for the basal branch in a set of 2638
327 local ML phylogenies, we found results that agree with the whole genome ASTRAL, SNAPP
328 and NJ trees: the most common basal branches are the pelagic groups *Rhamphochromis* and
329 *Diplotaxodon* (in 42.12% of the genomic windows). In comparison, *A. calliptera* (including IOC
330 samples) were found to be basal only in 5.99% of the windows (Supplementary Fig. 19).

331 Joyce et al.²⁰ reported that the mtDNA haplogroup of Indian Ocean catchment (IOC) *A.*
332 *calliptera* clustered with mbuna (as we confirm in Supplementary Fig. 15) and suggested that
333 there had been repeated colonization of Lake Malawi by two independent *Astatotilapia* lineages
334 with different mitochondrial haplogroups: the first founding the entire species flock, and the
335 second, with the IOC mtDNA haplogroup, introgressing into the Malawi radiation and
336 contributing strongly to the mbuna. This hypothesis predicts that, compared with the Malawi
337 catchment *A. calliptera*, the IOC *A. calliptera* should be closer to mbuna. However, across the
338 nuclear genome we found a strong signal in the opposite direction, with 30% excess allele
339 sharing between Malawi catchment *A. calliptera* and mbuna (Fig. 4d). Therefore, the Joyce et al.
340 hypothesis that the mbuna, the most species rich group within the radiation, may be a hybrid
341 lineage formed from independent invasions is not supported by genome-wide data.

342 It has been repeatedly suggested that *A. calliptera* may be the direct descendant of the riverine-
343 generalist lineage that seeded the Lake Malawi radiation^{7,50,53,54}. Our interpretation of this
344 argument is that the ancestor was likely a riverine generalist that was ecologically and
345 phenotypically similar to *A. calliptera* and other *Astatotilapia*. This hypothesis is lent further
346 support by geometric morphometric analysis. Using 17 homologous body shape landmarks we
347 established that, despite the relatively large genetic divergence, *A. calliptera* is nested within the
348 morphospace of the other more distantly related but ecologically similar *Astatotilapia* species
349 (Fig. 4a,e), and these together have a central position within the morphological space of the Lake
350 Malawi radiation (Fig. 4e and Supplementary Fig. 16).

351 To reconcile the nested phylogenetic position of *A. calliptera* with its generalist ‘prototype’
352 phenotype, we propose a model where the Lake Malawi species flock consists of three separate
353 radiations splitting off from the lineage leading to *A. calliptera*. The relationships between the
354 major groups supported by the ASTRAL, SNAPP and NJ methods suggest that the pelagic

355 radiation was seeded first, then the benthic + utaka, and finally the rock-dwelling mbuna, all in a
356 relatively quick succession, followed by subsequent gene flow as described above (Fig. 4f; the
357 pelagic vs. utaka + benthic branching order is swapped in SVDquartets tree in Supplementary
358 Fig. 9b). Applying our per-generation mutation rate to observed genomic divergences we
359 obtained mean divergence time estimates between these lineages between 460 thousand years
360 ago (ka) [95%CI: (350ka to 990ka)] and 390ka [95%CI: (300ka to 860ka)] (Fig. 4f), assuming
361 three years per generation as in ref. 55. The point estimates all fall within the second most recent
362 prolonged deep lake phase inferred from the Lake Malawi paleoecological record⁵⁶ while the
363 upper ends of the confidence intervals cover the third deep lake phase at ~800ky. Considering
364 that our split time estimates from sequence divergence are likely to be reduced by subsequent
365 gene-flow, leading to underestimates, the data are consistent with a previous report based on
366 fossil time calibration which put the origin of the Lake Malawi radiation at 700-800ka¹².

367 The fact that the common ancestor of all the *A. calliptera* appears younger than the Malawi
368 radiation suggests that the Lake Malawi *A. calliptera* population has been a reservoir that has
369 repopulated the river systems and more transient lakes following dry-wet transitions in East
370 African hydroclimate^{56,57}. Our results do not fully resolve whether the lineage leading from the
371 common ancestor to *A. calliptera* retained its riverine generalist phenotype throughout or
372 whether a lacustrine species evolved at some point (e.g. the common ancestor of *A. calliptera*
373 and mbuna) and later de-specialized again to recolonize the rivers. However, while it is a
374 possibility, we suggest it is unlikely that the many strong phenotypic affinities of *A. calliptera* to
375 the basal *Astatotilapia* (Fig. 4e; refs. ^{58,59}) would be reinvented from a lacustrine species.

376

377 **Signatures and consequences of selection on coding sequences**

378 To gain insight into the functional basis of diversification and adaptation in Lake Malawi
379 cichlids, we next turned our attention to protein coding genes. We compared the between-species
380 levels of non-synonymous variation \bar{p}_N to synonymous variation \bar{p}_S in 20,664 genes and
381 calculated the difference between these two values ($\delta_{N-S} = \bar{p}_N - \bar{p}_S$). Overall, coding sequence
382 exhibits signatures of purifying selection: the average between-species \bar{p}_N was 54% lower than
383 in a random matching set of non-coding regions. Interestingly, the average between-species
384 synonymous variation \bar{p}_S in genes was 13% higher than in non-coding control regions ($p <$
385 2.2×10^{-16} , one tailed Mann-Whitney test). One possible explanation of this observation would

386 be if intergenic regions were homogenized by gene-flow, whereas protein coding genes were
387 more resistant to this.

388 To control for statistical effects of variation in gene length and sequence composition we
389 normalized the δ_{N-S} values per gene by taking into account the variance across all pairwise
390 sequence comparisons for each gene, deriving the non-synonymous excess score Δ_{N-S} (see
391 Methods). Values at the upper tail of the distribution of Δ_{N-S} are substantially overrepresented in
392 the actual data when compared to a null model based on random sampling of codons (Fig. 5a).
393 We focus below on the top 5% of the distribution ($\Delta_{N-S} > 40.2$, 1034 candidate genes). Genes
394 with elevated Δ_{N-S} are expected to have been under positive selection at multiple non-
395 synonymous sites, either recently repeatedly within multiple species or ancestrally. Therefore,
396 the statistic reveals only a limited subset of positive selection events from the history of the
397 radiation (e.g. a selection event on a single amino acid would not be detected). Furthermore, to
398 minimise any effect of gene prediction errors, all the following analyses focus on the 15980
399 (77.3% of total) genes for which zebrafish homologs were found in ref. 11; selection scores of
400 genes without homologs are briefly discussed in a Supplementary Note.

401 Cichlids have an unexpectedly large number of gene duplicates, which possibly contributed to
402 their extensive adaptive radiations^{3,11}. To investigate the extent of divergent selection on gene
403 duplicates, we examined how the Δ_{N-S} scores are related to gene copy numbers in the reference
404 genomes. Focusing on homologous genes annotated both in the Malawi reference (*M. zebra*) and
405 in the zebrafish genome, we found that the highest proportion of candidate genes was among
406 genes with two or more copies in both genomes (N - N). The relative enrichment in this category
407 is both substantial and highly significant (Fig. 5b). On the other hand, the increase in proportion
408 of candidate genes in the N - 1 category (multiple copies in the *M. zebra* genome but only one
409 copy in zebrafish) is of a much lesser magnitude and is not significant (χ^2 test $p = 0.18$),
410 suggesting that selection is occurring more often within ancient multi-copy gene families, rather
411 than on genes with cichlid-specific duplications.

412 We used Gene Ontology (GO) annotation of zebrafish homologs to test whether candidate genes
413 are enriched for particular functional categories (Methods). We found significant enrichment for
414 30 GO terms (range: $1.6 \times 10^{-8} < p < 0.01$, `weigh` algorithm⁶⁰; Supplementary Table 3): 10 in the

415 Molecular Function, 4 in the Cellular Component and 16 in Biological Process category.
416 Combining all the results in a network (connecting terms that share many genes) revealed clear
417 clusters of enriched terms related to (i) haemoglobin function and oxygen transport; (ii)
418 phototransduction and visual perception; and (iii) the immune system, especially inflammatory
419 response and cytokine activity (Fig. 5c). That evolution of genes in these functional categories
420 has contributed to cichlid radiations has been suggested previously (see below); it is nevertheless
421 interesting that these categories stand out in a genome-wide analysis.

422

423 **Shared mechanisms of depth adaptation**

424 To gain insight into the distribution of adaptive alleles across the radiation, we built maximum
425 likelihood trees from amino acid sequences of candidate genes, thus summarising potentially
426 complex haplotype genealogy networks. Focusing on the significantly enriched GO categories,
427 many haplotype trees have features that are unusual in the broader dataset: the haplotypes from
428 the deep benthic group and the deep-water pelagic *Diplotaxodon* tend to group together (despite
429 these two groups being distant in whole-genome phylogenies and monophyletic in only two out
430 of 2638 local ML trees) and also tend to be disproportionately diverse when compared with the
431 rest of the radiation. We quantified both excess similarity and diversity, and found that both
432 measures are elevated for candidate genes in the ‘visual perception’ category (Fig. 6a; Mann-
433 Whitney tests: $p=0.007$ for similarity, $p=0.08$ for shared diversity, and $p=0.003$ when the scores
434 are added) and also for the ‘haemoglobin complex’ category (p values not significant due to the
435 small number of genes).

436 Sharply decreasing levels of dissolved oxygen and low light intensities with narrow short
437 wavelength spectra are the hallmarks of the habitats at below ~50 meters to which the deep
438 benthic and *Diplotaxodon* groups have both adapted, either convergently or in parallel⁶¹. Shared
439 signatures of selection in genes involved in vision and in oxygen transport therefore point to
440 shared molecular mechanisms underlying this ecological parallelism. Further evidence of shared
441 mechanisms of adaptation is that, for genes annotated with ‘photoreceptor activity’ and
442 ‘haemoglobin complex’ GO terms, the Δ_{N-S} selection score is strongly correlated with the local
443 levels of excess allele sharing between the two depth-adapted groups (Fig. 6b; $\rho_S = 0.63$ and
444 0.81 , $p = 0.001$ and $p = 0.051$, respectively).

445 Vision genes with high similarity and diversity scores for the deep benthic and *Diplotaxodon*
446 groups include three opsins: the green sensitive RH2A β and RH2B, and rhodopsin (Fig. 6a and
447 Supplementary Fig. 20). The specific residues that distinguish the deep adapted groups from the
448 rest of the radiation differ between the two RH2 copies, with only one shared mutation out of a
449 possible fourteen (Supplementary Fig. 20). RH2A β and RH2B are located less than 40kb apart
450 on the same chromosome (Fig. 6c); a third paralog, RH2A α , is located between them, but does
451 not show signatures of shared depth adaptation (Supplementary Fig. 21), consistent with reports
452 of functional divergence between RH2A α and RH2A β ^{62,63}. A similar, albeit weaker signature of
453 shared depth-related selection is apparent in rhodopsin, which is known to play a role in deep-
454 water adaptation in cichlids⁶⁴. Previously, we discussed the role of coding variants in rhodopsin
455 in the early stages of speciation of *A. calliptera* in the crater Lake Massoko⁵⁵. The haplotype tree
456 presented here for the broader radiation shows that the Massoko alleles did not originate by
457 mutation in that lake but were selected out of ancestral variation (Fig. 6a). The remaining opsin
458 genes are less likely to be involved in shared depth adaptation (Supplementary Note).

459 There have been many studies of selection on opsin genes in fish⁶⁵⁻⁶⁷, including selection
460 associated with depth preference, but having whole genome coverage allows us to investigate
461 other components of primary visual perception in an unbiased fashion. We found shared patterns
462 of selection between deep benthics and *Diplotaxodon* in six other vision associated candidate
463 genes (Fig. 6a). The functions of these genes, together with the fact that RH2A β and RH2B are
464 expressed exclusively in double-cone photoreceptors, suggest a prominent role of cone cell
465 vision in depth adaptation. The wavelength of maximum absorbance in cone cells expressing a
466 mixture of RH2A β with RH2B ($\lambda_{\max} = 498\text{nm}$) corresponds to the part of light spectrum that
467 transmits the best into deep water in Lake Malawi⁶⁷.

468 Figure 6c illustrates interactions of the vision genes with shared selection patterns in the cichlid
469 double-cone photoreceptor. The homeobox protein *six7* governs the expression of RH2 opsins
470 and is essential for the development of green cones in zebrafish⁶⁸ (specific mutations are
471 highlighted in Supplementary Fig. 20). The kinase GRK7 and the retinal cone arrestin-C have
472 complementary roles in photoresponse recovery: arrestin produces the final shutoff of the cone
473 pigment following phosphorylation by GRK7, thus determining the temporal resolution of
474 motion vision⁶⁹. Bases near to the C-terminus in RH2A β mutated away from serine (S290Y and

475 S292G), thus reducing the number of residues that can be modified by GRK7 (Supplementary
476 Fig. 20). The transducin subunit GNAT2 is located exclusively in the cone receptors and is a key
477 component of the pathway which converts light stimulus into electrical response in these cells⁷⁰.
478 Finally, peripherin-2 is essential to the development and renewal of the membrane system that
479 holds the opsin pigments in both rod and cone cells⁷¹.

480 Haemoglobin genes in teleost fish are found in two separate chromosomal locations: the minor
481 'LA' cluster and the major 'MN' cluster⁷². The region around the LA cluster has been
482 highlighted by selection scans among four *Diplotaxodon* species by Hahn et al.⁷³, who also noted
483 the similarity of the haemoglobin subunit beta (HB β) haplotypes between *Diplotaxodon* and deep
484 benthic species. We confirmed signatures of selection in the two annotated LA cluster
485 haemoglobins. In addition, we found that four haemoglobin subunits (HB β 1, HB β 2, HB α 2,
486 HB α 3) from the MN cluster are also among the genes with high selection scores (Supplementary
487 Fig. 22). The shared patterns of depth selection may be particular to the β -globin genes
488 (Supplementary Fig. 22), although this hypothesis remains tentative, because the repetitive
489 nature of the MN cluster precludes us from confidently examining all haemoglobin genes.

490 A key question concerns the mechanism leading to the similarity of haplotypes in *Diplotaxodon*
491 and deep benthics. Possibilities include parallel selection on variation segregating in both groups
492 due to common ancestry, selection on the gene flow that we described in a previous section, or
493 independent selection on new mutations. From considering the haplotype trees and local patterns
494 of excess allele sharing (f_{dM} statistics⁵⁵), there is evidence for each of these processes acting on
495 different genes. The haplotype trees for rhodopsin and HB β have outgroup taxa (and also *A.*
496 *calliptera*) appearing at multiple locations on their haplotype networks (Fig. 6a), suggesting that
497 the haplotype diversity of these genes may reflect ancestral variation. In contrast, trees for the
498 green cone genes show the Malawi radiation all being derived with respect to outgroups and we
499 found substantially elevated f_{dM} scores extending for around 40kb around the RH2 cluster (Fig.
500 6d), consistent with adaptive introgression in a pattern reminiscent of mimicry loci in *Heliconius*
501 butterflies⁷⁴. Finally, the peaks in f_{dM} around peripherin-2 and one of the arrestin-C genes are
502 narrow, ending at the gene boundaries, and f_{dM} scores are elevated only for non-synonymous
503 variants; synonymous variants do not show excess allele sharing (Supplementary Fig. 23). Due
504 to the close proximity of non-synonymous and synonymous sites within the same gene, this

505 suggests that for these two genes there may have been independent selection on the same *de novo*
506 mutations.

507 508 **Discussion**

509 Variation in genome sequences forms the substrate for evolution. Here we described genome
510 variation at the full sequence level across the Lake Malawi haplochromine cichlid radiation. We
511 focused on ecomorphological diversity, representing more than half the genera from each major
512 group, rather than obtaining deep coverage of species within any particular group. Therefore, we
513 have more samples from the morphologically highly diverse benthic lineages than, for example,
514 the mbuna where there are relatively fewer genera and many species are largely recognised by
515 colour differences.

516 The observation that cichlids within an African Great Lake radiation are genetically very similar
517 is not new⁷⁵, but we now quantify the relationship of this to within-species variation, and the
518 consequences for variation in local phylogeny across the genome. The fact that between-species
519 divergence is generally only slightly higher than within species diversity, is likely the result of
520 the young age of the radiation, the relatively low mutation rate, and of gene flow between taxa.
521 Within-species diversity itself is relatively low for vertebrates, at around 0.1%, suggesting that
522 low genome-wide nucleotide diversity levels do not necessarily limit rapid adaptation and
523 speciation. This conclusion appears in contrast for example with a recent report that high
524 diversity levels may have been important for rapid adaptation in Atlantic killifish⁷⁶. One
525 possibility is that in cichlids repeated selection has maintained diversity in adaptive alleles for a
526 range of traits that support ecological diversification, as we have concluded for rhodopsin and
527 HB β and appears to be the case for some adaptive variants in sticklebacks⁷⁷.

528 We provide evidence that gene flow during the radiation, although not ubiquitous, has certainly
529 been extensive. Overall, the numerous violations of the bifurcating species tree model suggest
530 that full resolution of interspecies relationships in this system will require network approaches
531 (see e.g. ref. 35; section 6.2) and population genomic analyses within the framework of the
532 structured coalescent with gene flow. The majority of the signals affect groups of species,
533 suggesting events involving their common ancestors, or are between closely-related species
534 within the major ecological groups. The only strong and clear example of recent gene flow

535 between individual distantly-related species is not within Lake Malawi itself, but between
536 *Otopharynx tetrastigma* from crater Lake Ilamba and local *A. calliptera*. Lake Ilamba is very
537 turbid and the scenario is reminiscent of cichlid admixture in low visibility conditions in Lake
538 Victoria⁷⁸. It is possible that some of the earlier signals of gene flow between lineages we
539 observed in Lake Malawi may have happened during low lake level periods when the water is
540 known to have been more turbid⁵⁶.

541 Our model of the early stages of radiation in Lake Malawi (Fig. 4f) is broadly consistent with the
542 model of initial separation by major habitat divergence²³, although we propose a refinement in
543 which there were three relatively closely-spaced separations from a generalist *Astatotilapia* type
544 lineage, initially of pelagic genera *Rhamphochromis* and *Diplotaxodon*, then of shallow- and
545 deep-water benthics and utaka (this includes Kocher's sand dwellers^{23,29}), and finally of mbuna.
546 Thus, we suggest that Lake Malawi contains three separate haplochromine cichlid radiations
547 stemming from the generalist lineage, interconnected by subsequent gene flow.

548 The finding that cichlid-specific gene duplicates do not tend to diverge particularly strongly in
549 coding sequences (Fig. 5b) suggests that other mechanisms of diversification following gene
550 duplications may be more important. Divergence via changes in expression patterns has
551 previously been illustrated and discussed¹¹, and future studies addressing structural variation
552 between cichlid genomes will assess the contribution of differential retention of duplicated
553 genes.

554 The evidence concerning shared adaptation of the visual and oxygen-transport systems to deep-
555 water environments between deep benthics and *Diplotaxodon* suggests different evolutionary
556 mechanisms acting on different genes, even within the same cellular system. It will be interesting
557 to see whether the same genes or even specific mutations underlie depth adaptation in Lake
558 Tanganyika, which harbours specialist deep water species in least two different tribes⁷⁹ and has a
559 similar light attenuation profile but a steeper oxygen gradient than Lake Malawi⁶¹.

560 Over the last few decades, East African cichlids have emerged as a model for studying rapid
561 vertebrate evolution^{11,23}. Taking advantage of recently assembled reference genomes¹¹, our data
562 and results provide unprecedented information about patterns of sequence sharing and adaptation
563 across the Lake Malawi radiation, with insights into mechanisms of rapid phenotypic

564 diversification. The data sets are openly available (see Acknowledgements) and will underpin
565 further studies on specific taxa and molecular systems. For example, we envisage that our
566 results, clarifying the relationships between all the main lineages and many individual species,
567 will facilitate speciation studies, which require investigation of taxon pairs at varying stages on
568 the speciation continuum^{80,81}, and studies on the role of adaptive gene flow in speciation.

569 **Methods**

570 **Samples.** Ethanol preserved fin clips were collected by M.J. Genner and G.F. Turner between
571 2004 and 2014 from Tanzania and Malawi, in collaboration with the Tanzania Fisheries
572 Research Institute (the MolEcoFish Project) and with the Fisheries Research Unit of the
573 Government of Malawi (various collaborative projects). Samples were collected and exported
574 with the permission of the Tanzania Commission for Science and Technology, the Tanzania
575 Fisheries Research Institute, and the Fisheries Research Unit of the Government of Malawi

576 **From sequencing to a variant callset.** The analyses presented above are based on SNPs
577 obtained from Illumina short (100bp-125bp) reads, aligned to the *Metriaclima zebra* reference
578 assembly version 1.1¹¹ with `bwa-mem`⁸², followed by GATK haplotype caller⁸³ and
579 `samtools/bcftools`⁸⁴ variant calling restricted to 653Mb of ‘accessible genome’ where variants
580 can be determined confidently with short reads, filtering, genotype refinement, imputation, and
581 phasing in BEAGLE⁸⁵ and further haplotype phasing with `shapeit v2`⁸⁶, including the use of
582 phase-informative reads⁸⁷. For details please see Supplementary Methods.

583
584 **Linkage disequilibrium (LD) calculations.** Haplotype r^2 between pairs of SNPs was calculated
585 along the phased scaffolds 0 to 201, using `vcftools v0.1.12b` with the options `--hap-`
586 `r2 --ld-window-bp 50000`. To reduce the computational burden, we used a random
587 subsample of 10% of SNPs. We binned the r^2 values according to the distance between SNPs
588 into 1kb or 100bp windows and plotted the average values in each bin.

589 To estimate background LD, we calculated haplotype r^2 between variants mapping to
590 different linkage groups (LG) in the *Oreochromis niloticus* genome assembly. First, we used the
591 chain files generated by the whole genome alignment pipeline⁸⁸ (see Supplementary Methods)
592 and the UCSC `liftOver` tool to translate the genomic coordinates of all SNPs to the *O.*
593 *niloticus* coordinates. Then we calculated LD between variants mapping to LG1 and LG2.

594
595 **De novo mutation rate estimation.** In each trio we looked for mutations in the child that were
596 not present in either of its parents. Because the results of this analysis are very sensitive to false
597 positives and false negative rates, we used higher coverage sequencing (~40x average) and
598 applied more stringent genome masks than in the population genomic work. Increased coverage
599 supports clean separation of sequencing errors and somatic mutations from true heterozygous
600 calls in the offspring, and improved ability to distinguish single copy vs. multi-copy sequence on
601 a per-individual basis.

602 First we determined the “Accessible Genome” (i.e. the regions of the genome where we
603 can confidently call *de novo* mutations) for each trio by excluding:

- 604 1. Genomic regions where mapped read depth in any member of a trio is $\leq 25\times$ or $>50\times$

- 605 2. Bases where either of the parents has a mapped read that does not match the reference
606 (the specific bases where any read has non-reference alleles in the parents were masked)
607 3. Sequences where indels were called in any sample (we also excluded +/- 3bp of sequence
608 surrounding the indel)
609 4. Sites which were called as multiallelic among the nine samples in the overall trios dataset
610 5. Known segregating variable sites - i.e. sites with alternative alleles found in four and
611 more copies in the overall Lake Malawi dataset
612 6. Sites in the reference where less than 90% of overlapping 50-mers (sub-sequences of
613 length 50) could be matched back uniquely and without 1-difference. For this we used
614 Heng Li's SNPable tool (<http://lh3lh3.users.sourceforge.net/snpable.shtml>), dividing the
615 reference genome into overlapping k-mers (sequences of length k – we used k=50), and
616 then aligning the extracted k-mers back to the genome (we used `bwa aln -R`
617 `1000000 -O 3 -E 3`).

618 After excluding sites in the categories above, we were left with an “Accessible Genome” of
619 516.6Mb in the *A. calliptera* trio, 459Mb in the *A. stuartgranti* trio and 404Mb in the *L. lethrinus*
620 trio. Because any observed *de novo* mutation could have occurred either on the chromosome
621 inherited from the mother or on the chromosome inherited from the father, the point estimate of
622 the per generation per basepair mutation rate is: $\mu = n\text{Mutations}/(2 \times \text{AccessibleGenome})$.

623 Next we set out to search for *de novo* mutations: i.e. heterozygous sites in the offspring
624 within the Accessible Genome. Under random sampling there is an equal probability of seeing a
625 read with either of the two alleles at a heterozygous site. Therefore, N_a - the number of reads
626 supporting the alternative allele is distributed as: $\sim \text{Binomial}(\text{ReadDepth}, 0.5)$. We filtered out
627 variants with observed N_a below 2.5th or above 97.5th percentiles of this distribution, thus
628 accepting a false negative rate of 5%. We also filtered out sites where the offspring call had Read
629 Position or Base Quality rank sum test Z-score > 99.5 th percentile of standard normal
630 distribution or where the Strand Bias phred scaled p-value was ≥ 20 or where the phred scaled
631 GQ (genotype quality) in either mother, father, or offspring was ≤ 30 . For simplicity, assuming
632 these filters are independent they are expected to introduce a false negative rate of 7.17%. The
633 mutation rate estimate was adjusted to account for this.

634 After filtering, we found nine *de novo* mutations across the three offspring. For each
635 mutation we double-checked the alignment in the IGV genome browser and found all of them
636 were single base mutations supported by high number of reads (>8) in the offspring. The 95%
637 confidence intervals for the number of observed mutations were calculated using the “exact”
638 method relating chi-squared and Poisson distributions^{89,90}. If N is the number of observed
639 mutations, the lower (ciN_L) and upper (ciN_U) limits are:

$$640 \quad ciN_L = \frac{P(\chi^2_{2N} \leq 0.025)}{2} \quad ciN_U = \frac{P(\chi^2_{2(N+1)} \geq 0.975)}{2}$$

641 where $2N$ and $2(N+1)$ are degrees of freedom of the corresponding chi-squared distributions.

642
643 **Principal Component Analysis.** SNPs with minor allele frequency ≥ 0.05 were selected using
644 the `bcftools (v1.2) view option --min-af 0.05:minor`. The program `vcftools`
645 `v0.1.12b` was then used to export that data into PLINK format⁹¹. Next, the variants were LD-
646 pruned to obtain a set of variants in approximate linkage equilibrium (unlinked sites) using the
647 `--indep-pairwise 50 5 0.2` option in PLINK v1.0.7. Principal Component Analysis
648 on the resulting set of variants was performed using the `smartpca` program from the
649 `eigensoft v5.0.2` software package⁹² with default parameters.

650 **Genome-wide F_{ST} calculations.** In addition to performing PCA, the `smartpca` program from
651 the `eigensoft v5.0.2` software package also calculates genome-wide F_{ST} for all pairs of
652 populations specified by the sixth column in the `.pedind` file. For the calculation, it uses the
653 Hudson estimator, as defined by Bhatia, Patterson et al.⁹³ in equation 10, and the ‘ratio of
654 averages’ is used to combine estimates of F_{ST} across multiple variants, as recommended in that
655 manuscript. We used all SNPs (no minor allele frequency filtering).

656 **Allele sharing test for group assignment.** We tested if two individuals who come from the
657 same group always share more derived alleles with each other than with any individuals from
658 other groups. Technically, we implemented this using the D statistic (ABBA-BABA tests)
659 framework^{31,32}, by calculating $D(A, G1, G2, O)$ for all permutations of individuals, where $G1$ and
660 $G2$ come from the same eco-morphological group and A from a different group. The outgroup O
661 was always *N. brichardi* from Lake Tanganyika. Note that this is an unusual use of the D
662 statistics and our aim here was not to look for gene flow but to test if allele sharing is greater
663 within eco-morphological groups ($G1$ with $G2$) compared to across groups (A with $G2$), in
664 which case $D(A, G1, G2, O) > 0$. All results were statistically significant, which was assessed
665 using block jackknife³¹ on windows of 60k SNPs .

666 **D_{min} statistic.** Here we calculated the D statistic for each trio of species (A,B,C) and for all
667 possible tree topologies (the outgroup again fixed as *N. brichardi*). Therefore,
668 $D_{min} = \min(|D(A,B,C,O)|, |D(A,C,B,O)|, |D(C,B,A,O)|)$. If this is significantly elevated, then allele
669 sharing within the trio of species is inconsistent with any simple tree topology. Note that this
670 approach is conservative in the sense that the D_{min} score for each trio is considered in isolation
671 and we ignore ‘higher-order’ inconsistencies where different D_{min} trio topologies are inconsistent
672 with each other. Statistical significance was assessed using block jackknife³¹ on windows of 60k
673 SNPs and FWER was calculated following the Holm-Bonferroni method.

674 **Sample selection for demographic analyses.** To prevent potential confounding effects of
675 uneven sequencing depth, we limited these analyses to one high coverage (15x) individual per
676 species. Species without a high coverage sample (*P. subocularis*, *F. rostratus* and *L. trewavasae*)
677 were not included.

678
679 **Outgroup sequences/alleles.** Outgroup (Supplementary Table 5) sequences in *M. zebra* genomic
680 coordinates were obtained based on pairwise whole-genome alignments (Supplementary
681 Methods). Insertions in the outgroup were ignored and deletions filled by N characters.

682
683 **Local phylogenetic trees and maximum clade credibility.** To generate a multiple alignment
684 input in `fasta` format we used the `getWGSeq` subprogram of `evo`. We set the window size in
685 terms of the numbers of variants rather than physical length (8000 variants; `--split 8000`
686 option) aiming for the local regions to have similar strengths of phylogenetic signal. Small
687 windows at the ends of scaffolds were discarded. We limited the sequence output to the
688 accessible genome using the `--accessibleGenomeBED` option. The *N. brichardi* outgroup
689 sequence in *M. zebra* genomic coordinates was added via the `--incl-Pn` option.

690 Maximum likelihood phylogenies were inferred using RAxML v7.7.8⁹⁴ under the
691 GTRGAMMA model. The best tree for each region was selected out of twenty alternative runs
692 on distinct starting maximum parsimony trees (the -N 20 option).

693 The maximum clade credibility (MCC) trees were calculated in TreeAnnotator
694 v.2.4.2, a part of the BEAST2 platform⁹⁵. Clade credibility is the frequency with which a
695 clade appears in the tree set; the MCC tree is the tree (from among the trees in the set) that
696 maximizes the product of the frequencies of all its clades³³. The node heights for the MCC trees
697 are derived as a summary from the heights of each clade in the whole tree set via the Common
698 Ancestor heights option.

699
700 **Mitochondrial (mtDNA) phylogenies.** The mtDNA sequence corresponds to scaffolds 747 and
701 2036 in the *Metriaclima zebra* reference. Variants from these scaffolds were subjected to the
702 same filtering as in the rest of the genome except for the depth filter because the mapped read
703 depth was much higher (approximately 300-400x per sample). Because of the greater sequence
704 diversity in the mtDNA genome, we found that more than 10% of variants were multiallelic.
705 Therefore, we separated SNPs from indels at multiallelic sites using `bcftools norm` with the
706 `--multiallelics -` option, then removed indels and the merged multiallelic SNPs back
707 together with the `--multiallelics +` option. Sequences in the `fasta` format were
708 generated using the `bcftools consensus` command, and missing genotypes in the VCF
709 replaced by the N character with the `--mask` option. *N. brichardi* outgroup sequence in *M.*
710 *zebra* genomic coordinates was added to the `fasta` files.

711 A maximum likelihood tree was inferred using RAxML v7.7.8⁹⁴ under the
712 GTRGAMMA model. The best tree was selected out of twenty alternative runs on distinct
713 starting maximum parsimony trees (using the -N 20 option) and two hundred bootstrap
714 replicates were obtained using RAxML's rapid bootstrapping algorithm⁹⁶ satisfying the -N
715 `autoFC` frequency-based bootstrap stopping criterion. Bipartition bootstrap support was drawn
716 on the maximum likelihood tree using the RAxML `-f b` option.

717
718 **Neighbour-joining trees and the residuals.** For the Neighbour-Joining (NJ)⁹⁷ trees we
719 calculated the average numbers of single nucleotide differences between haplotypes for each pair
720 of species. This simple pairwise difference matrix was divided by the accessible genome size to
721 obtain pairwise differences per bp, which are equivalent to \hat{p}_{AB} of Dasarathy et al.⁴². Then we
722 followed equation 8 from Dasarathy et al. and calculated their *corrected* measure of
723 dissimilarity:

$$\hat{d}_{AB} = -\frac{3}{4} \log \left(1 - \frac{4}{3} \hat{p}_{AB} \right)$$

724 The \hat{d}_{AB} values were then used as input into the `nj()` tree-building function implemented in the
725 APE package⁹⁸ in R language.

726 We measured the distances between all pairs of species in the reconstructed NJ tree (i.e.
727 the lengths of branches) using the `get_distance()` method implemented in the ETE3 toolkit
728 for phylogenetic trees⁹⁹. Our first measure of 'tree violation' is the difference between these
729 distances and the distances between samples in the original matrix that was used to build the NJ
730 tree.

731

732 **Multispecies coalescent methods.** We applied three different methods that attempt to
733 reconstruct the species tree under the multispecies coalescent model. For a brief discussion of
734 these approaches see Supplementary Methods.

735 For SNAPP³⁶ we used a random subset of ~0.5% of genome-wide SNPs (48,922 SNPs)
736 for 12 individuals representing the eco-morphological groups and the Lake Victoria outgroup *P.*
737 *nyererei* whose alleles were filled in based on the whole genome alignment. The *P. nyererei*
738 alleles were assigned as ‘ancestral’ (0 in the nexus input file). The ‘forward’ and ‘backward’
739 mutation rate parameters u and v were calculated directly from the data by SNAPP (the `Calc`
740 `mutation rates` option). The default value 10 was used for the Coalescent rate parameter
741 and the value of the parameter was sampled (estimated in the MCMC chain). We used
742 uninformative priors as we don’t assume strong a priori knowledge about the parameters. The
743 prior for ancestral population sizes was chosen to be a relatively broad gamma distribution with
744 parameters $\alpha = 4$ and $\beta = 20$. The tree height prior λ was set to the initial value of 100 but
745 sampled in the MCMC chain with an uninformative uniform hyperprior on the interval
746 [0,50000]. We ran three independent MCMC chains with the same starting parameters, each on
747 30 threads with a total runtime of over 10 CPU years. The first one million steps from each
748 MCMC chain was discarded as burn-in. In total, more than 30 million MCMC steps were
749 sampled in the three runs. For the MCMC traces for each run see Supplementary Fig. 24.

750 Next we used SVDquartets^{37,38} as implemented in PAUP* [v4.0a (build 159)]¹⁰⁰. We
751 prepared the data into the NEXUS ‘dna’ format, using `evo` with the `getWGSeq --whole-`
752 `genome --makeSVDinput -r` options. This command outputs for each individual the
753 DNA base at each variable site, randomly sampling one of the two alleles at heterozygous sites,
754 and ignoring sites that become monomorphic due to this random sampling of alleles. The final
755 dataset contained 17,833,187 SNPs. Then we ran SVDquartets in PAUP* setting outgroup to *N.*
756 *brichardi* and then executing `svdq evalq=all`; specifying that all quartets should be
757 evaluated (not just a random subset). In the final step, PAUP* version of the QFM algorithm¹⁰¹ is
758 used to search for the overall tree that minimizes the number of quartets that are inconsistent
759 with it.

760 Finally we used ASTRAL⁴⁰ (v. 5.6.1) with default parameters and the full set of 2543
761 local trees generated by RAxML (see above) as input.

762
763 **Tree comparisons.** To summarise the degree of (dis)agreement between the topologies of trees
764 produced by different phylogenetic methods (Fig. 2c), we calculated the normalised Robinson-
765 Foulds distances between pairs of trees¹⁰² using the `RF.dist` function from the `phangorn`¹⁰³
766 package in R with the option `normalize=TRUE`.

767
768 **Chromopainter and fineSTRUCTURE.** Singleton SNPs were excluded using the `bcftools`
769 `v.1.1 -c 2:minor` option, before exporting the remaining variants in the PLINK format⁹¹.
770 The `chromopainter v0.0.4` software⁴⁷ was then run for the 201 largest genomic scaffolds
771 on `shapeit` phased SNPs. Briefly, we created a uniform recombination map using the
772 `makeuniformrecfile.pl` script, then estimated the effective population size (N_e) for a
773 subsample of 20 individuals using the `chromopainter` inbuilt expectation-maximization
774 procedure⁴⁷, averaged over the 20 N_e values using the provided `neaverage.pl` script. The
775 `chromopainter` program was then run for each scaffold independently, with the `-a 0 0`
776 option to run all individuals against all others. Results for individual scaffolds were combined

777 using the chromocombine tool before running fineSTRUCTURE v0.0.5 with 1,000,000
778 burn in iterations, and 200,000 sample iterations, recording a sample every 1,000 iterations
779 (options -x 1000000 -y 200000 -z 1000). Finally, the sample relationship tree was
780 built with fineSTRUCTURE using the -m T option and 20,000 iterations.

781
782 **The f -branch statistic.** The f_4 -admixture ratio (f statistic) statistic was developed to estimate the
783 proportion of introgressed material in an admixed population [see SOM18 in ref. 31, and f_G in
784 ref. 48]. However, when calculated for different subsets of samples within the same phylogeny,
785 there are a very large number of highly correlated f values that are hard to interpret. To make the
786 interpretation easier, we developed “ f -branch” or $f_b(C)$: $f_b(C) = \text{median}_A[\min_B[f(A, B, C, O)]]$,
787 where B are samples descending from branch b , and A are samples descending from the sister
788 branch of b . The outgroup O was always *N. brichardi*. The $f_b(C)$ score provides for each branch b
789 of a given phylogeny and each sample C a summary of excess allele sharing of branch b with
790 sample C (Fig 3, Supplementary Figure 26). Each $f_b(C)$ score was also assigned an associated z -
791 score to assess statistical significance $Z_b(C) = \text{median}_A[\min_B[Z(A, B, C, O)]]$. Additional
792 information on the f and $f_b(C)$ statistics, including detailed reasoning behind the design of $f_b(C)$,
793 are in Supplementary Methods.

794
795 **Geometric morphometric analyses.** A total of 168 photographs were used to compare the gross
796 body morphology of *Astatotilapia calliptera* to that of endemic Lake Malawi species and other
797 East African *Astatotilapia* lineages (Supplementary Table 7). Coordinates for 17 homologous
798 landmarks [following ref. 104] were collected using tpsDig2 v2.26¹⁰⁵. After landmark
799 digitization, analysis of shape variation was carried out in R (v3.3.2) using the package
800 GeOMorph v3.0.2¹⁰⁶. First a General Procrustes Analysis was applied to remove non-shape
801 variation and shape data were corrected for allometric size effects by performing a regression of
802 Procrustes coordinates (10,000 iterations). The resulting allometry corrected residuals were used
803 in PCA.

804
805 **Maps.** Present day catchment boundary maps are based on ‘level 3’ detail of Hydro1K dataset
806 from the US Geological Survey. We downloaded the watershed boundary data from the United
807 Nations Environment Programme website (<http://ede.grid.unep.ch>) and processed it using the
808 QGIS geographic information system software (<http://www.qgis.org/en/site/>).

809
810 **Protein-coding gene annotations.** We used the BROADMZ2 annotation generated by the cichlid
811 genome project¹¹ and removed overlapping transcripts using Jim Kent’s
812 genePredSingleCover program. Genes whose annotated length in nucleotides was not
813 divisible by three were discarded, as they typically had inaccuracies in annotation that would
814 require manual curation (2495 out of 23698 genes). We also used the cichlid genome project¹¹
815 assignment of homologs between the *M. zebra* genome reference and zebrafish (*Danio rerio*).

816
817 **Coding sequence positive selection scan.** We used evo with the getCodingSeq -H b --
818 no-stats options to obtain the coding sequences for each allele and each gene. The excess of
819 non-synonymous variation (δ_{N-S}) and the non-synonymous variation excess score (Δ_{N-S}) were
820 calculated on a per-gene basis as follows. Let N_{TS} be the number of possible non-synonymous
821 transitions and N_{TV} the number of possible non-synonymous transversion between two

822 sequences; analogously S_{TS} and S_{TV} represent possible synonymous differences. We do not
 823 specify the ancestral allele, and therefore consider it equally likely that allele i mutated into allele
 824 j or that allele j mutated to allele i . Then let N_d be the number of observed non-synonymous
 825 mutations and S_d the number of observed synonymous mutations. If there are more than one
 826 difference within a codon, all “mutation pathways” (i.e. the different orders in which mutations
 827 could have happened) have equal probabilities. When a particular allele contained a premature
 828 stop codon, the remainder of the sequence after the stop was excluded from the calculations.

829 Because the transition:transversion ratio in the Lake Malawi dataset was 1.73, and hence
 830 (because there are two possible transversions for each possible transition) the prior probability of
 831 each transition is 3.46 times that of each transversion, we account for the unequal probabilities of
 832 transitions and transversions in calculating the proportions of non-synonymous (p_N) and of
 833 synonymous differences (p_S) as follows:

$$p_N = \frac{N_d}{3.46 \times N_{TS} + N_{TV}} \quad p_S = \frac{S_d}{3.46 \times S_{TS} + S_{TV}}$$

834
 835 The excess of non-synonymous variation (δ_{N-S}) is the average of $p_N - p_S$ over pairwise
 836 sequence comparisons. Only between-species sequence comparisons are considered for the Lake
 837 Malawi dataset. We normalized the δ_{N-S} values in order to take into account the effect on the
 838 variance of this statistic introduced by differences in gene length and by sequence composition.
 839 To achieve this, we used the leave-one-out jackknife procedure across different pairwise
 840 comparisons for each gene, estimating the standard error. The non-synonymous variation excess
 841 score (Δ_{N-S}) is then:

$$\Delta_{N-S} = \frac{\delta_{N-S}}{jackknife_se(\delta_{N-S})}$$

842
 843 Note that because the sequences are related by a genealogy, there is a correlation structure
 844 between the pairwise comparisons. Therefore, the jackknife approach substantially
 845 underestimates the true standard error of δ_{N-S} and is used here simply as a normalization factor.

846 The null model shown in Fig. 5a was derived by splitting all the coding sequence into its
 847 constituent codons, and then randomly sampling these codons with replacement to build new
 848 sequences that matched the actual coding genes in their numbers and the length distribution.
 849 Then we calculated the Δ_{N-S} scores, as we did for the actual genes and compared the two
 850 distributions. High positive values at the upper tail of the distribution are substantially
 851 overrepresented in the actual data when compared to a null model.

852 We also calculated the above statistics for random non-coding regions, matching the gene
 853 sequences in length. We used the `bedtools v2.26.0107 shuffle` command to permute the
 854 locations of exons along the chromosomes. Of the total length of all the permuted sequences,
 855 98.4% was within the ‘accessible genome’ and outside of coding sequences (we required at least
 856 95% in any of the permuted locations). The specific command was:

```
857 bedtools shuffle -chrom -I exons.bed -excl InaccessibleGenome_andExons.bed -f  
858 0.05 -g chrom.sizes
```

859
 860 **Gene Ontology enrichment.** Zebrafish has the most extensive functional gene annotation of any
 861 fish species, providing a basis for Gene Ontology (GO)¹⁰⁸ term enrichment analysis. Gene
 862 Ontology (GO) enrichment for the genes that were candidates for being under positive selection
 863 (the top 5% of Δ_{N-S} values) was calculated in R using the `topGO v2.26.0` package¹⁰⁹ from

864 the Bioconductor project¹¹⁰. The GO hierarchical structure was obtained from the GO.db
865 v3.4.0 annotation and linking zebrafish gene identifiers to GO terms was accomplished using
866 the org.Dr.eg.db v3.4.0 annotation package. Genome-wide, between 9024 and 9353
867 genes had a GO annotation that could be used by topGO, the exact number depending on the
868 GO category being assessed. The nodeSize parameter was set to 5 to remove GO terms which
869 have fewer than 5 annotated genes, as suggested in the topGO manual.

870 There is often an overlap between gene sets annotated with different GO terms, in part
871 because the terms are related to each other in a hierarchical structure¹⁰⁸. This is partly accounted
872 for by our use in topGO of the weight algorithm that accounts for the GO graph structure by
873 down-weighting genes in the GO terms that are neighbours of the locally most significant terms
874 in the GO graph⁶⁰. All the p-values we report are from the weight algorithm, which the authors
875 suggest should be reported without multiple testing correction¹⁰⁹.

876 Some interdependency between significant GO terms remains after using the weight
877 algorithm. Therefore, we used the Enrichment Map¹¹¹ app for Cytoscape
878 (<http://www.cytoscape.org>) to organize all the significantly enriched terms into networks where
879 terms are connected if they have a high overlap, i.e. if they share many genes.

880

881 ***Diplotaxodon* and deep benthic convergence.** To obtain a quantitative measure of the similarity
882 between and the extent of excess diversity in the *Diplotaxodon* and deep benthic amino acid
883 sequences, we calculated simple statistics based on the proportions of non-synonymous
884 differences (p_N scores). Intuitively, the similarity score is high if *Diplotaxodon* and deep benthic
885 jointly have higher p_N than all the others, but are not very different from each other relative to
886 how much diversity there is within *Diplotaxodon* and deep benthic.

887 Specifically, the similarity score s is calculated as follows:

$$s_{raw} = \bar{p}_N^O - (\bar{p}_N^B - \bar{p}_N^W)$$

888 and

$$s = \frac{s_{raw}}{jackknife_se(p_N)} - mean\left(\frac{s_{raw}}{jackknife_se(p_N)}\right)$$

889 where \bar{p}_N^O is the mean p_N between *Diplotaxodon* jointly with deep benthic and all the other Lake
890 Malawi species, \bar{p}_N^B is the mean p_N between *Diplotaxodon* and deep benthic, and \bar{p}_N^W is the mean
891 p_N within *Diplotaxodon* and deep benthic. The jackknife normalization is analogous to the one
892 used for Δ_{N-S} and the mean (\bar{s}_{raw}) is subtracted to center the statistic at zero.

893 The excess diversity score is high when the mean p_N scores within *Diplotaxodon* and
894 within deep benthic are high relative to the mean p_N in the rest of the radiation. Specifically, the
895 excess score ex is defined as:

$$ex = \frac{[(\bar{p}_N^D + \bar{p}_N^{DB})/2] - \bar{p}_N^R}{jackknife_se(p_N)}$$

896

897 where \bar{p}_N^D is the mean p_N within *Diplotaxodon*, \bar{p}_N^{DB} is the mean p_N within deep benthic, and \bar{p}_N^R
898 is the mean p_N within the rest of the radiation.

899

900 **Haplotype trees.** To view the relationship between haplotypes for genes of interest, we
901 translated nucleotide sequences to amino acid sequences and loaded these into Haplotype
902 Viewer (<http://www.cibiv.at/~greg/haploviewer>). This software requires that a tree is loaded

903 together with the sequences. Therefore, we inferred gene trees using RAxML v7.7.8⁹⁴ with the
904 PROTGAMMADAYHOFF model of substitution.

905
906 **Local excess allele sharing between *Diplotaxodon* and deep benthic.** We used an extension of
907 the f_d statistic⁴⁸; this extension is referred to as f_{dM} ⁵⁵. f_{dM} is a conservative version of the f statistic
908 that is particularly suited for analysis of small genomic windows^{48,55}. For the gene scores shown
909 in Fig. 6b, we calculated f_{dM} (mbuna, deep benthic, *Diplotaxodon*, *N. brichardi*) for each gene in
910 window from the transcription start site (TSS) to 10 kb into the gene. For the ‘along the genome’
911 plots, as shown in Fig. 6d and Supplementary Fig. 23, we used a product of two f_{dM} statistics
912 [f_{dM} (shallow benthic, deep benthic, *Diplotaxodon*, *N. brichardi*) x f_{dM} (*Rhamphochromis*,
913 *Diplotaxodon*, deep benthic, *N. brichardi*)], an approach which we found to increase the local
914 resolution. This score was calculated in sliding windows of 100 SNPs across a region of +
915 100kb around the genes. Finally, we also calculated f_{dM} (mbuna, deep benthic, *Diplotaxodon*, *N.*
916 *brichardi*) separately for synonymous and non-synonymous mutations in each gene.

917
918 **Reporting Summary.** Further information on experimental design is available in the Nature
919 Research Reporting Summary linked to this article.

920
921 **Code availability.** The majority of the custom code used in this project is available on Github as
922 a part of the evo package (<https://github.com/millanek/evo>). All other custom codes are available
923 from the authors upon request.

924
925 **Data availability.** All raw sequencing reads have been deposited to the NCBI Short Read
926 Archive: (BioProjects PRJEB1254 and PRJEB15289). Sample accessions are listed in
927 Supplementary Table 4. In addition, we are making whole-genome variant calls in the Variant
928 Call Format (VCF), phylogenetic trees and protein coding sequence alignments, and tables with
929 f_4 statistics available through the Dryad Digital Repository
930 (<http://dx.doi.org/10.5061/dryad.7rj8k6c>).

931

932 **References and Notes:**

- 933 1. Losos, J. B. & Ricklefs, R. E. Adaptation and diversification on islands. *Nature* **457**,
934 830–836 (2009).
- 935 2. Wagner, C. E., Harmon, L. J. & Seehausen, O. Ecological opportunity and sexual
936 selection together predict adaptive radiation. *Nature* **487**, 366–369 (2012).
- 937 3. Berner, D. & Salzburger, W. The genomics of organismal diversification illuminated by
938 adaptive radiations. *Trends Genet.* **31**, 491–499 (2015).
- 939 4. Darwin, C. *On the Origin of Species*. (OUP Oxford, 2008).
- 940 5. Lamichhaney, S. *et al.* Evolution of Darwin's finches and their beaks revealed by
941 genome sequencing. *Nature* **518**, 371–375 (2015).
- 942 6. Losos, J., Jackman, T., Larson, A., Queiroz, K. & Rodriguez-Schettino, L. Contingency
943 and determinism in replicated adaptive radiations of island lizards. *Science* **279**, 2115–
944 2118 (1998).
- 945 7. Fryer, G. & Iles, T. D. *The cichlid fishes of the great lakes of Africa: their biology and*
946 *evolution*. (Oliver and Boyd, 1972).
- 947 8. Salzburger, W., Van Bocxlaer, B. & Cohen, A. S. Ecology and Evolution of the African

- 948 Great Lakes and Their Faunas. *Annu. Rev. Ecol. Evol. Syst.* **45**, 519–545 (2014).
- 949 9. Genner, M. J. *et al.* How does the taxonomic status of allopatric populations influence
950 species richness within African cichlid fish assemblages? *Journal of Biogeography* **31**,
951 93–102 (2004).
- 952 10. Meyer, A. Phylogenetic relationships and evolutionary processes in East African cichlid
953 fishes. *Trends Ecol Evol* **8**, 279–284 (1993).
- 954 11. Brawand, D. *et al.* The genomic substrate for adaptive radiation in African cichlid fish.
955 *Nature* **513**, 375–381 (2014).
- 956 12. Meyer, B. S., Matschiner, M. & Salzburger, W. Disentangling incomplete lineage
957 sorting and introgression to refine species-tree estimates for Lake Tanganyika cichlid
958 fishes. *Syst. Biol.* (2016). doi:10.1093/sysbio/syw069
- 959 13. Meier, J. I. *et al.* Ancient hybridization fuels rapid cichlid fish adaptive radiations. *Nat*
960 *Commun* **8**, 14363 (2017).
- 961 14. Koblmüller, S., Egger, B., Sturmbauer, C. & Sefc, K. M. Rapid radiation, ancient
962 incomplete lineage sorting and ancient hybridization in the endemic Lake Tanganyika
963 cichlid tribe Tropheini. *Mol. Phylogenet. Evol.* **55**, 318–334 (2010).
- 964 15. Weiss, J. D., Cotterill, F. P. D. & Schliewen, U. K. Lake Tanganyika—A ‘Melting Pot’
965 of Ancient and Young Cichlid Lineages (Teleostei: Cichlidae)? *PLoS ONE* **10**,
966 e0125043 (2015).
- 967 16. Gante, H. F. *et al.* Genomics of speciation and introgression in Princess cichlid fishes
968 from Lake Tanganyika. *Mol Ecol* (2016). doi:10.1111/mec.13767
- 969 17. Wagner, C. E. *et al.* Genome-wide RAD sequence data provide unprecedented
970 resolution of species boundaries and relationships in the Lake Victoria cichlid adaptive
971 radiation. *Mol Ecol* **22**, 787–798 (2012).
- 972 18. Genner, M. J. & Turner, G. F. Ancient hybridization and phenotypic novelty within
973 Lake Malawi's cichlid fish radiation. *Mol. Biol. Evol.* **29**, 195–206 (2012).
- 974 19. Moran, P., Kornfield, I. & Reinthal, P. N. Molecular Systematics and Radiation of the
975 Haplochromine Cichlids (Teleostei: Perciformes) of Lake Malawi. *Copeia* **1994**, 274
976 (1994).
- 977 20. Joyce, D. A. *et al.* Repeated colonization and hybridization in Lake Malawi cichlids. **21**,
978 R108–9 (2011).
- 979 21. Leffler, E. M. *et al.* Revisiting an old riddle: what determines genetic diversity levels
980 within species? *PLoS Biol.* **10**, e1001388 (2012).
- 981 22. Albertson, R. C., Markert, J. A., Danley, P. D. & Kocher, T. D. Phylogeny of a rapidly
982 evolving clade: the cichlid fishes of Lake Malawi, East Africa. *Proc. Natl. Acad. Sci.*
983 *U.S.A.* **96**, 5107–5110 (1999).
- 984 23. Kocher, T. D. Adaptive evolution and explosive speciation: the cichlid fish model. *Nat.*
985 *Rev. Genet.* **5**, 288–298 (2004).
- 986 24. Chimpanzee Sequencing and Analysis Consortium. Initial sequence of the chimpanzee
987 genome and comparison with the human genome. *Nature* **437**, 69–87 (2005).
- 988 25. Ravi, V. & Venkatesh, B. Rapidly evolving fish genomes and teleost diversity. *Curr.*
989 *Opin. Genet. Dev.* **18**, 544–550 (2008).
- 990 26. Recknagel, H., Elmer, K. R. & Meyer, A. A hybrid genetic linkage map of two
991 ecologically and morphologically divergent Midas cichlid fishes (*Amphilophus* spp.)
992 obtained by massively parallel DNA sequencing (ddRADSeq). *G3 (Bethesda)* **3**, 65–74
993 (2013).

- 994 27. Ségurel, L., Wyman, M. J. & Przeworski, M. Determinants of mutation rate variation in
995 the human germline. *Annu Rev Genomics Hum Genet* **15**, 47–70 (2014).
- 996 28. Feng, C. *et al.* Moderate nucleotide diversity in the Atlantic herring is associated with a
997 low mutation rate. *Elife* **6**, e23907 (2017).
- 998 29. Konings, A. *Malawi Cichlids in Their Natural Habitat*. (Cichlid Press, 2007).
- 999 30. 1000 Genomes Project Consortium *et al.* An integrated map of genetic variation from
1000 1,092 human genomes. *Nature* **491**, 56–65 (2012).
- 1001 31. Green, R. E. *et al.* A draft sequence of the Neandertal genome. *Science* **328**, 710–722
1002 (2010).
- 1003 32. Durand, E. Y., Patterson, N., Reich, D. & Slatkin, M. Testing for ancient admixture
1004 between closely related populations. *Mol. Biol. Evol.* **28**, 2239–2252 (2011).
- 1005 33. Heled, J. & Bouckaert, R. R. Looking for trees in the forest: summary tree from
1006 posterior samples. *BMC Evol. Biol.* **13**, 221 (2013).
- 1007 34. Edwards, S. V. Is a new and general theory of molecular systematics emerging?
1008 *Evolution* **63**, 1–19 (2009).
- 1009 35. Edwards, S. V. *et al.* Implementing and testing the multispecies coalescent model: A
1010 valuable paradigm for phylogenomics. *Mol. Phylogenet. Evol.* **94**, 447–462 (2016).
- 1011 36. Bryant, D., Bouckaert, R., Felsenstein, J., Rosenberg, N. A. & RoyChoudhury, A.
1012 Inferring species trees directly from biallelic genetic markers: bypassing gene trees in a
1013 full coalescent analysis. *Mol. Biol. Evol.* **29**, 1917–1932 (2012).
- 1014 37. Chifman, J. & Kubatko, L. Quartet Inference from SNP Data Under the Coalescent
1015 Model. *Bioinformatics* **30**, 3317–3324 (2014).
- 1016 38. Chifman, J. & Kubatko, L. Identifiability of the unrooted species tree topology under the
1017 coalescent model with time-reversible substitution processes, site-specific rate variation,
1018 and invariable sites. *J. Theor. Biol.* **374**, 35–47 (2015).
- 1019 39. Long, C. & Kubatko, L. The Effect of Gene Flow on Coalescent-Based Species-Tree
1020 Inference. *arXiv.org* 1–19 (2017).
- 1021 40. Mirarab, S. *et al.* ASTRAL: genome-scale coalescent-based species tree estimation.
1022 *Bioinformatics* **30**, i541–8 (2014).
- 1023 41. Zhang, C., Rabiee, M., Sayyari, E. & Mirarab, S. ASTRAL-III: polynomial time species
1024 tree reconstruction from partially resolved gene trees. *BMC Bioinformatics* **19**, 153
1025 (2018).
- 1026 42. Dasarathy, G., Nowak, R. & Roch, S. Data Requirement for Phylogenetic Inference
1027 from Multiple Loci: A New Distance Method. *IEEE/ACM Trans Comput Biol Bioinform*
1028 **12**, 422–432 (2015).
- 1029 43. Rusinko, J. & McPartlon, M. Species tree estimation using Neighbor Joining. *J. Theor.*
1030 *Biol.* **414**, 5–7 (2017).
- 1031 44. Ballard, J. W. O. & Whitlock, M. C. The incomplete natural history of mitochondria.
1032 *Mol Ecol* **13**, 729–744 (2004).
- 1033 45. Toews, D. P. L. & Brelsford, A. The biogeography of mitochondrial and nuclear
1034 discordance in animals. *Mol Ecol* **21**, 3907–3930 (2012).
- 1035 46. Consuegra, S., John, E., Verspoor, E. & de Leaniz, C. G. Patterns of natural selection
1036 acting on the mitochondrial genome of a locally adapted fish species. *Genet. Sel. Evol.*
1037 **47**, 58 (2015).
- 1038 47. Lawson, D. J., Hellenthal, G., Myers, S. & Falush, D. Inference of Population Structure
1039 using Dense Haplotype Data. *PLoS Genet.* **8**, e1002453 (2012).

- 1040 48. Martin, S. H., Davey, J. W. & Jiggins, C. D. Evaluating the use of ABBA-BABA
1041 statistics to locate introgressed loci. *Mol. Biol. Evol.* **32**, 244–257 (2015).
- 1042 49. Martin, S. H. *et al.* Genome-wide evidence for speciation with gene flow in Heliconius
1043 butterflies. *Genome Res.* **23**, 1817–1828 (2013).
- 1044 50. Eccles, D. H. & Trewavas, E. *Malawian Cichlid Fishes*. (Lake Fish Movies, 1989).
- 1045 51. Eriksson, A. & Manica, A. Effect of ancient population structure on the degree of
1046 polymorphism shared between modern human populations and ancient hominins. *Proc.*
1047 *Natl. Acad. Sci. U.S.A.* **109**, 13956–13960 (2012).
- 1048 52. Pickrell, J. K. & Pritchard, J. K. Inference of population splits and mixtures from
1049 genome-wide allele frequency data. *PLoS Genet.* **8**, e1002967 (2012).
- 1050 53. Peterson, E. N., Cline, M. E., Moore, E. C., Roberts, N. B. & Roberts, R. B. Genetic sex
1051 determination in *Astatotilapia calliptera*, a prototype species for the Lake Malawi cichlid
1052 radiation. *Naturwissenschaften* **104**, 41 (2017).
- 1053 54. Genner, M. J., Ngatunga, B. P., Mzighani, S., Smith, A. & Turner, G. F. Geographical
1054 ancestry of Lake Malawi's cichlid fish diversity. *Biol. Lett.* **11**, 20150232 (2015).
- 1055 55. Malinsky, M. *et al.* Genomic islands of speciation separate cichlid ecomorphs in an East
1056 African crater lake. *Science* **350**, 1493–1498 (2015).
- 1057 56. Ivory, S. J. *et al.* Environmental change explains cichlid adaptive radiation at Lake
1058 Malawi over the past 1.2 million years. *Proceedings of the National Academy of*
1059 *Sciences* **113**, 11895–11900 (2016).
- 1060 57. Lyons, R. P. *et al.* Continuous 1.3-million-year record of East African hydroclimate, and
1061 implications for patterns of evolution and biodiversity. *Proc. Natl. Acad. Sci. U.S.A.* **112**,
1062 15568–15573 (2015).
- 1063 58. Greenwood, P. H. *Towards a phyletic classification of the 'genus' Haplochromis*
1064 *(Pisces, Cichlidae) and related taxa. Part 1.* **35**, 265–322 (1979).
- 1065 59. Lippitsch, E. A phyletic study on lacustrine haplochromine fishes (Perciformes,
1066 Cichlidae) of East Africa, based on scale and squamation characters. *Journal of Fish*
1067 *Biology* **42**, 903–946 (1993).
- 1068 60. Alexa, A., Rahnenführer, J. & Lengauer, T. Improved scoring of functional groups from
1069 gene expression data by decorrelating GO graph structure. *Bioinformatics* **22**, 1600–
1070 1607 (2006).
- 1071 61. Van Bocxlaer, B., SCHULTHEIß, R., Plisnier, P.-D. & Albrecht, C. Does the decline of
1072 gastropods in deep water herald ecosystem change in Lakes Malawi and Tanganyika?
1073 *Freshwater Biology* **57**, 1733–1744 (2012).
- 1074 62. Spady, T. C. *et al.* Evolution of the cichlid visual palette through ontogenetic
1075 subfunctionalization of the opsin gene arrays. *Mol. Biol. Evol.* **23**, 1538–1547 (2006).
- 1076 63. Weadick, C. J. & Chang, B. S. W. Complex patterns of divergence among green-
1077 sensitive (RH2a) African cichlid opsins revealed by Clade model analyses. *BMC Evol.*
1078 *Biol.* **12**, 206 (2012).
- 1079 64. Sugawara, T. *et al.* Parallelism of amino acid changes at the RH1 affecting spectral
1080 sensitivity among deep-water cichlids from Lakes Tanganyika and Malawi. *Proc. Natl.*
1081 *Acad. Sci. U.S.A.* **102**, 5448–5453 (2005).
- 1082 65. Bowmaker, J. K. & Hunt, D. M. Evolution of vertebrate visual pigments. *Current*
1083 *Biology* **16**, R484–R489 (2006).
- 1084 66. Davies, W. I. L., Collin, S. P. & Hunt, D. M. Molecular ecology and adaptation of visual
1085 photopigments in craniates. *Mol Ecol* **21**, 3121–3158 (2012).

- 1086 67. Carleton, K. L., Dalton, B. E., Escobar-Camacho, D. & Nandamuri, S. P. Proximate and
1087 ultimate causes of variable visual sensitivities: Insights from cichlid fish radiations.
1088 *Genesis* **54**, 299–325 (2016).
- 1089 68. Ogawa, Y., Shiraki, T., Kojima, D. & Fukada, Y. Homeobox transcription factor Six7
1090 governs expression of green opsin genes in zebrafish. *Proceedings of the Royal Society*
1091 *B: Biological Sciences* **282**, 20150659 (2015).
- 1092 69. Renninger, S. L., Gesemann, M. & Neuhauss, S. C. F. Cone arrestin confers cone vision
1093 of high temporal resolution in zebrafish larvae. *Eur. J. Neurosci.* **33**, 658–667 (2011).
- 1094 70. Brockerhoff, S. E. *et al.* Light stimulates a transducin-independent increase of
1095 cytoplasmic Ca²⁺ and suppression of current in cones from the zebrafish mutant *nof*. *J.*
1096 *Neurosci.* **23**, 470–480 (2003).
- 1097 71. Boesze-Battaglia, K. & Goldberg, A. F. X. Photoreceptor renewal: a role for
1098 peripherin/rds. *Int. Rev. Cytol.* **217**, 183–225 (2002).
- 1099 72. Opazo, J. C., Butts, G. T., Nery, M. F., Storz, J. F. & Hoffmann, F. G. Whole-genome
1100 duplication and the functional diversification of teleost fish hemoglobins. *Mol. Biol.*
1101 *Evol.* **30**, 140–153 (2013).
- 1102 73. Hahn, C., Genner, M. J., Turner, G. F. & Joyce, D. A. The genomic basis of cichlid fish
1103 adaptation within the deepwater ‘twilight zone’ of Lake Malawi. *Evolution Letters* **1**,
1104 184–198 (2017).
- 1105 74. Heliconius Genome Consortium. Butterfly genome reveals promiscuous exchange of
1106 mimicry adaptations among species. *Nature* **487**, 94–98 (2012).
- 1107 75. Meyer, A., Kocher, T. D., Basasibwaki, P. & Wilson, A. C. Monophyletic origin of Lake
1108 Victoria cichlid fishes suggested by mitochondrial DNA sequences. *Nature* **347**, 550–
1109 553 (1990).
- 1110 76. Reid, N. M. *et al.* The genomic landscape of rapid repeated evolutionary adaptation to
1111 toxic pollution in wild fish. *Science* **354**, 1305–1308 (2016).
- 1112 77. Jones, F. C. *et al.* The genomic basis of adaptive evolution in threespine sticklebacks.
1113 *Nature* **484**, 55–61 (2012).
- 1114 78. Seehausen, O. Cichlid Fish Diversity Threatened by Eutrophication That Curbs Sexual
1115 Selection. *Science* **277**, 1808–1811 (1997).
- 1116 79. Konings, A. *Tanganyika cichlids in their natural habitat*. (Cichlid Press, 2015).
- 1117 80. Coyne, J. A. & Orr, H. A. Speciation. (2004).
- 1118 81. Feder, J. L., Egan, S. P. & Nosil, P. The genomics of speciation-with-gene-flow. *Trends*
1119 *in Genetics* **28**, 342–350 (2012).
- 1120 82. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-
1121 MEM. *arXiv.org q-bio.GN*, (2013).
- 1122 83. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for
1123 analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
- 1124 84. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping
1125 and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**,
1126 2987–2993 (2011).
- 1127 85. Browning, S. R. & Browning, B. L. Rapid and accurate haplotype phasing and missing-
1128 data inference for whole-genome association studies by use of localized haplotype
1129 clustering. *Am. J. Hum. Genet.* **81**, 1084–1097 (2007).
- 1130 86. Delaneau, O., Marchini, J. & Zagury, J.-F. A linear complexity phasing method for
1131 thousands of genomes. *Nat. Methods* **9**, 179–181 (2012).

- 1132 87. Delaneau, O., Howie, B., Cox, A. J., Zagury, J.-F. & Marchini, J. Haplotype estimation
1133 using sequencing reads. *Am. J. Hum. Genet.* **93**, 687–696 (2013).
- 1134 88. Miller, W. *et al.* 28-way vertebrate alignment and conservation track in the UCSC
1135 Genome Browser. *Genome Res.* **17**, 1797–1808 (2007).
- 1136 89. Ulm, K. A simple method to calculate the confidence interval of a standardized mortality
1137 ratio (SMR). *Am. J. Epidemiol.* **131**, 373–375 (1990).
- 1138 90. Dobson, A. J., Kuulasmaa, K., Eberle, E. & Scherer, J. Confidence intervals for
1139 weighted sums of Poisson parameters. *Stat Med* **10**, 457–462 (1991).
- 1140 91. Purcell, S. *et al.* PLINK: A Tool Set for Whole-Genome Association and Population-
1141 Based Linkage Analyses. *The American Journal of Human Genetics* **81**, 559–575
1142 (2007).
- 1143 92. Patterson, N., Price, A. L. & Reich, D. Population structure and eigenanalysis. *PLoS*
1144 *Genet.* **2**, e190–e190 (2006).
- 1145 93. Bhatia, G., Patterson, N., Sankararaman, S. & Price, A. L. Estimating and interpreting
1146 FST: The impact of rare variants. *Genome Res.* **23**, 1514–1521 (2013).
- 1147 94. Stamatakis, A. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses
1148 with thousands of taxa and mixed models. *Bioinformatics* **22**, 2688–2690 (2006).
- 1149 95. Bouckaert, R. *et al.* BEAST 2: a software platform for Bayesian evolutionary analysis.
1150 *PLoS Comput. Biol.* **10**, e1003537 (2014).
- 1151 96. Stamatakis, A., Hoover, P. & Rougemont, J. A rapid bootstrap algorithm for the RAxML
1152 Web servers. *Syst. Biol.* **57**, 758–771 (2008).
- 1153 97. Saitou, N. & Nei, M. The neighbor-joining method: a new method for reconstructing
1154 phylogenetic trees. *Mol. Biol. Evol.* **4**, 406–425 (1987).
- 1155 98. Paradis, E., Claude, J. & Strimmer, K. APE: Analyses of Phylogenetics and Evolution in
1156 R language. *Bioinformatics* **20**, 289–290 (2004).
- 1157 99. Huerta-Cepas, J., Serra, F. & Bork, P. ETE 3: Reconstruction, Analysis, and
1158 Visualization of Phylogenomic Data. *Mol. Biol. Evol.* **33**, 1635–1638 (2016).
- 1159 100. Swofford, D. L. *PAUP*: Phylogenetic analysis using parsimony (and other methods)*,
1160 2002. (Sinauer Associates).
- 1161 101. Reaz, R., Bayzid, M. S. & Rahman, M. S. Accurate Phylogenetic Tree Reconstruction
1162 from Quartets: A Heuristic Approach. *PLoS ONE* **9**, e104008
- 1163 102. Robinson, D. F. & Foulds, L. R. Comparison of phylogenetic trees. *Mathematical*
1164 *Biosciences* (1981).
- 1165 103. Schliep, K. P. phangorn: phylogenetic analysis in R. *Bioinformatics* **27**, 592–593 (2011).
- 1166 104. Theis, A., Ronco, F., Indermaur, A., Salzburger, W. & Egger, B. Adaptive divergence
1167 between lake and stream populations of an East African cichlid fish. *Mol Ecol* **23**, 5304–
1168 5322 (2014).
- 1169 105. Rohlf, F. J. tpsDig2.
- 1170 106. Adams, D. C. & Castillo, E. O. geomorph: an R package for the collection and analysis
1171 of geometric morphometric shape data. *Methods in Ecology and ...* (2013).
- 1172 107. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing
1173 genomic features. *Bioinformatics* **26**, 841–842 (2010).
- 1174 108. Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. The Gene
1175 Ontology Consortium. *Nat. Genet.* **25**, 25–29 (2000).
- 1176 109. Alexa, A. & Rahnenfuhrer, J. topGO: enrichment analysis for gene ontology. *R package*
1177 *version* (2010).

- 1178 110. Huber, W. *et al.* Orchestrating high-throughput genomic analysis with Bioconductor.
1179 *Nat. Methods* **12**, 115–121 (2015).
1180 111. Merico, D., Isserlin, R., Stueker, O., Emili, A. & Bader, G. D. Enrichment map: a
1181 network-based method for gene-set enrichment visualization and interpretation. *PLoS*
1182 *ONE* **5**, e13984 (2010).
1183
1184
1185
1186
1187
1188

1189
1190
1191
1192
1193
1194
1195
1196
1197

Acknowledgments

This work was supported by the Wellcome Trust (097677/Z/11/Z to MM, WT206194 and WT207492 for RD and HS), the Royal Society-Leverhulme Trust Africa Awards (AA100023 and AA130107 to MG and GFT) and the European Molecular Biology Organization (ALTF 456-2016 to MM). We want to thank the Sanger Institute sequencing core for DNA sequencing, Mingliu Du for DNA extractions, David Swofford and Michael Matschiner for advice on phylogenomic analyses, Walter Salzburger and Ian Wilson for comments on the manuscript. We also thank the Tanzania Fisheries Research Institute, and the Fisheries Research Unit of the Government of Malawi for their assistance and support.

1198
1199
1200
1201
1202

Author contributions:

EM, GFT, MJG, MM and RD devised the study. GFT and MJG collected the samples. AMT bred parent-offspring trios and performed geometric morphometric analyses. MM performed the DNA extractions. HS and MM analysed the genomic data. All authors participated in interpretation of the results. MM, HS and RD drafted the manuscript, and all others commented.

1203
1204
1205

Competing interests

RD declares that he owns stock in Illumina from previous consulting. The authors declare no other competing interests.

1206 **Figure captions:**

1207
1208 **Fig. 1: The Lake Malawi haplochromine cichlid radiation.** **a**, The sampling coverage of this study: overall and
1209 for each of the seven main eco-morphological groups within the radiation. A representative specimen is shown for
1210 each group (*Diplotaxodon*: *D. limnothrissa*; shallow benthic: *Lethrinops albus*; deep benthic: *Lethrinops gossei*;
1211 mbuna: *Metriaclima zebra*; utaka: *Copadichromis virginalis*; *Rhamphochromis*: *R. woodi*). Numbers of species and
1212 genera are based on ref. 29. **b**, The distributions of genomic sequence diversity within individuals (heterozygosity;
1213 π) and of divergence between species (d_{XY}). **c**, Principal component analysis (PCA) of whole genome variation data.

1214
1215 **Fig. 2: Excess allele sharing and patterns of species relatedness.** **a**, Derived allele sharing reveals non-tree-like
1216 relationships among trios of species. The bars show the proportion of significantly elevated D_{\min} scores (see main
1217 text). Shading corresponds to FWER q values of (from light to dark) 10^{-2} , 10^{-4} , 10^{-8} , 10^{-14} . The scatterplots show the
1218 D_{\min} scores that were significant at $\text{FWER} < 0.01$. Results are shown separately for comparisons where all three
1219 species in the trio are from the same group, and for cases where the species come from two or three different groups.
1220 *Rhamphochromis* and utaka within-group comparisons are not shown due to the low number of data points. **b**, A set
1221 of 2543 Maximum Likelihood (ML) phylogenetic trees for non-overlapping regions along the genome. Branch
1222 lengths were scaled for visualization so that the total height of each tree is the same. The local trees were built with
1223 71 species and then subsampled for display to 12 individuals representing the eco-morphological groups. The
1224 maximum clade credibility tree shown here was built from the subsampled local trees. A ML mitochondrial
1225 phylogeny is shown for comparison. **c**, A summary of all phylogenies from this study and the normalised Robinson-
1226 Foulds distances between them, reflecting the topological distance between pairs of trees on the scale from zero to
1227 100%. The least controversial 12 sample tree is SNAPP.t1, with an average distance to other trees of 17.7%, while
1228 ASTRAL* is the least controversial among the ‘main trees’ (mean distance of 25.3%). To compare trees with
1229 differing sets of taxa, the trees were downsampled so that only matching taxa were present. The position of the
1230 outgroup/root was considered in all comparisons.

1231 **Fig. 3: Identifying tree violating branches and possible gene flow events.** The branch-specific statistic $f_b(C)$
1232 identifies excess sharing of derived alleles between the branch of the tree on the y-axis and the species C on the x-
1233 axis (see Supplementary Note). The ASTRAL* tree was used as a basis for the branch statistic and grey data points
1234 in the matrix correspond to tests that are not consistent with the phylogeny. Colours correspond to eco-
1235 morphological groups as in Fig. 1. The * sign denotes block jack-knifing significance at $|Z| > 3.17$ (Holm-Bonferroni
1236 $\text{FWER} < 0.001$).

1237 **Fig. 4: Origins of the radiation and the role of *A. calliptera*.** **a**, An NJ phylogeny showing the Lake Malawi
1238 radiation in the context of other East African *Astatotilapia* taxa. **b**, A Lake Malawi NJ phylogeny with expanded
1239 view of *A. calliptera*, with all other groups collapsed. **c**, Approximate *A. calliptera* sampling locations shown on a
1240 map of the broader Lake Malawi region. Black lines correspond to present day level 3 catchment boundaries from
1241 the US Geological Survey’s HYDRO1k dataset. **d**, Strong f_4 admixture ratio signal showing that Malawi catchment

1242 *A. calliptera* are closer to mbuna than their Indian Ocean catchment counterparts. **e**, PCA of body shape variation of
1243 Lake Malawi endemics, *A. calliptera* and other *Astatotilapia* taxa, obtained from geometric morphometric analysis.
1244 **f**, A phylogeny with the same topology as in panel (b) but displayed with a straight line between the ancestor and *A.*
1245 *calliptera*. For each branch off this lineage, we show mean sequence divergence (d_{XY}) minus mean heterozygosity,
1246 and translation of this value into a mean divergence time estimate with 95% CI reflecting the statistical uncertainty
1247 in mutation rate. Dashed lines with arrows indicate likely instances of gene flow between major groups; their true
1248 timings are uncertain.

1249 **Fig. 5: Gene selection scores, copy numbers, and ontology enrichment.** **a**, The distribution of the non-
1250 synonymous variation excess scores (Δ_{N-S}) highlighting the top 5% cutoff, compared against a null model. The null
1251 was derived by calculating the statistic on randomly sampled combinations of codons. We also show the
1252 distributions of genes in selected Gene Ontology (GO) categories which are overrepresented in the top 5%. **b**, The
1253 relationship between the probability of Δ_{N-S} being in the top 5% and the relative copy numbers of genes in the Lake
1254 Malawi reference (*M. zebra*) and zebrafish. The p-values are based on χ^2 tests of independence. Genes existing in
1255 two or more copies in both zebrafish and Malawi cichlids are disproportionately represented among candidate
1256 selected genes. **c**, An enrichment map for significantly enriched GO terms (cutoff at $p \leq 0.01$). The level of overlap
1257 between GO enriched terms is indicated by the thickness of the edge between them. The colour of each node
1258 indicates the p-value for the term and the size of the node is proportional to the number of genes annotated with that
1259 GO category.

1260 **Fig. 6: Shared selection between the deep water adapted groups *Diplotaxodon* and deep benthic.** **a**, The
1261 scatterplot shows the distribution of genes with high Δ_{N-S} scores (candidates for positive selection) along axes
1262 reflecting shared selection signatures. Only genes with zebrafish homologs are shown. Amino acid haplotype trees,
1263 shown for genes as indicated by the red symbols and numbers, indicate that *Diplotaxodon* and deep benthic species
1264 are often divergent from other taxa, but similar to each other. Outgroups include *Oreochromis niloticus*,
1265 *Neolamprologus brichardi*, *Astatotilapia burtoni*, and *Pundamilia nyererei*. **b**, Selection scores plotted against f_{dM}
1266 (mbuna, deep benthic, *Diplotaxodon*, *N. brichardi*), a measure of local excess allele sharing between deep benthic
1267 and *Diplotaxodon*. Overall there is no correlation between Δ_{N-S} and f_{dM} . However, the strong correlation between
1268 Δ_{N-S} and f_{dM} in the highlighted GO categories suggests that positively selected alleles in those categories tend to be
1269 subject to introgression or convergent selection between *Diplotaxodon* and the deep benthic group. **c**, A schematic
1270 drawing of a double cone photoreceptor expressing the green sensitive opsins and illustrating the functions of other
1271 genes with signatures of shared selection. **d**, f_{dM} calculated in sliding windows of 100 SNPs around the green opsin
1272 cluster, revealing that excess allele sharing between deep benthic and *Diplotaxodon* extends far beyond the coding
1273 sequences.

1274
1275











