1 # Whole genome resequencing data enables a targeted

2 # SNP panel for conservation and aquaculture of

3 # *Oreochromis* cichlid fishes

4 Ciezarek A[1]*, Ford AGP[2]*, Etherington GJ[1], Kasozi N[3], Malinsky M[4], Mehta T[1], Penso-Dolfin L[5],

5 Ngatunga BP[6], Shechonge A[6], Tamatamah R[6], Haerty W[1], Di Palma F[7], Genner MJ[8], Turner

6 GF[9]

7    1. Earlham Institute, Norwich Research Park, Norwich, NR4 7UZ, UK.

8    2. Department of Life Sciences, Roehampton University, London, SW15 4JD, UK.

9    3. National Agricultural Research Organisation, Abi Zonal Agricultural Research and

10    Development Institute, P.O. Box 219, Arua, Uganda

11    4. Zoological Institute, Department of Environmental Sciences, University of Basel, 4051

12    Basel, Switzerland.

13    5. Silence Therapeutics GmbH, Robert-Rössle-Straße 10, 13125 Berlin, Germany.

14    6. Tanzania Fisheries Research Institute (TAFIRI) PO. Box 9750. Dar es Salaam. Tanzania.

15    7. School of Biological Sciences, University of East Anglia, Norwich, NR4 7TU, UK

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

17    9. School of Natural Sciences, Bangor University, Bangor, LL57 2UW, UK. Orchid: 0000-

18    0003-0099-7261

19    *Authors contributed equally

20    **Correspondence:**

21    adam.ciezarek@earlham.ac.uk, george.turner@bangor.ac.uk.

# Abstract

Cichlid fish of the genus *Oreochromis* form the basis of the global tilapia aquaculture and fisheries industry. Non-native farmed tilapia populations are known to be widely distributed across Africa and to hybridize with native *Oreochromis* species. However, many species are difficult to distinguish morphologically, hampering attempts to maintain good quality farmed strains or to identify pure populations of native species. Here, we describe the development of a single nucleotide polymorphism (SNP) genotyping panel from whole-genome resequencing data that enables targeted species identification in Tanzania. We demonstrate that an optimized panel of 96 genome-wide SNPs based on $F_{ST}$ outliers performs comparably to whole genome resequencing in distinguishing species and identifying hybrids. We also show this panel outperforms microsatellite-based and phenotype-based classification methods. Case studies indicate several locations where introduced aquaculture species have become established in the wild, threatening native *Oreochromis* species. The novel SNP markers identified here represent an important resource for assessing broodstock purity and helping to conserve unique endemic biodiversity, and in addition potentially for assessing broodstock purity in hatcheries.

**Keywords:** *Oreochromis*; tilapia; Tanzania; aquaculture; fisheries; hybridization; introduced species.

# Introduction

Global aquaculture production has increased rapidly in recent decades. Continued expansion is particularly important in Africa, where rapid human population growth over this century will stress food production systems (FAO, 2020). Tilapia, cichlid fish of the genus *Oreochromis*, native to Africa and the Middle East, have been a key part of the expansion of tropical aquaculture, accounting for 5.5 million of the global total of 47 million tonnes of inland finfish aquaculture production in 2018 (FAO, 2020). However, farmed populations have frequently colonized water catchments where they are not native, both due to deliberate introductions and accidental escape from fish farms (Shechonge, Ngatunga, Bradbeer, et al., 2019). This has threatened native species through ecological competition, habitat alteration and hybridization (Bbole et al., 2014; Canonico et al., 2005; Deines et al., 2014; Firmat et al., 2013; Macaranas et al., 1986; Ndiwa et al., 2014; Waiswa Mwanja et al., 2012). At present native *Oreochromis* species are poorly characterized, and their conservation could benefit from the identification of purebred populations for protection. Such safeguarding of the wild relatives of farmed species would also protect unique genetic resources that could be used to enhance traits in cultured *Oreochromis* strains (Macaranas et al., 1986; Thodesen et al., 2013).

Tanzania, a hotspot of natural diversity for tilapia species, has eight fully endemic *Oreochromis* species (*O. amphimelas*, *O. chungruruensis, O. karomo, O. korogwe, O. latilabris, O. ndalalani, O. rukwaensis, O. urolepis*). It also has an additional 12 species that are endemic to catchments shared with neighboring countries (*O. alcalicus, O. esculentus, O. girigan, O. hunteri, O. jipe, O. karongae, O. lidole, O. malagarasi, O. pangani, O. squamipinnis, O. tanganicae, O. variabilis*). Several of these species are adapted to unique environmental conditions, such as elevated temperatures, salinity, and pH (Ford et al., 2019; Trewavas, 1983). In addition, although Tanzania hosts a native population of *O. niloticus* indigenous to Lake

64    Tanganyika (Shechonge, Ngatunga, Tamatamah, et al., 2019), non-native farmed populations,

65    largely sourced from Lake Victoria, have been widely distributed across the country (Kajungiro et

66    al., 2019; Moses et al., 2020). The spread of *O. niloticus* has been accompanied by *O.*

67    *leucostictus*, another species present in Lake Victoria (Bradbeer et al., 2019; Shechonge,

68    Ngatunga, Bradbeer, et al., 2019; Shechonge et al., 2018). The Lake Victoria populations of

69    both *O. niloticus* and *O. leucostictus* were themselves introduced from the Nile system, mostly

70    likely Lake Albert, during the 1950s (Balirwa, 1988).

71         Nile tilapia (*O. niloticus*), in particular, is becoming established across Africa outside of its

72    natural range, including in South Africa (D'Amato et al., 2007), Zambia (Deines et al., 2014),

73    Zimbabwe (Marufu & Chifamba, 2013), the Democratic Republic of Congo (Goudswaard et al.,

74    2002; Mamonekene & Stiassny, 2012), Kenya (Angienda et al., 2011), as well as Tanzania

75    (Shechonge, Ngatunga, Bradbeer, et al., 2019), with reports of either replacement of the native

76    species or extensive introgressive hybridization (Bradbeer et al., 2019; Shechonge, Ngatunga,

77    Bradbeer, et al., 2019; Shechonge et al., 2018). There is also evidence of parasite transmission

78    from introduced tilapia species to native species (Jorissen et al., 2020). Despite this, intentional

79    movement and stocking of tilapia species into natural water bodies continues in many regions of

80    Africa (Genner et al., 2013).

81         Several studies have shown that diagnosis of *Oreochromis* hybrids purely based on

82    phenotypic traits of colour or morphology is unreliable (Bbole et al., 2014). Genetic analysis is

83    therefore necessary to determine if introduced and native strains are interbreeding, as well as

84    assessing broodstock purity in commercial aquaculture centres. Mitochondrial DNA has proved

85    insufficient for species diagnosis and resolution of many tilapia species due to recent

86    hybridization (Mojekwu et al. 2021). On the other hand, recent studies have shown the utility of

87    nuclear single nucleotide polymorphism (SNP) data for species and strain diagnosis between

88    species of *Oreochromis* (Syaifudin et al., 2019) and between strains of *O. niloticus* (Lind et al.,

89  2019). Meanwhile, high-throughput sequencing has proved useful in the development of

90  population-specific or species-diagnostic SNP panels for several commercially important

91  fisheries species, including Atlantic salmon (Campbell & Narum, 2011; Larson et al., 2014),

92  European herring (Helyar et al., 2012), Pacific lamprey (Hess et al., 2015) and white bass (Zhao

93  et al., 2019).

94      Here, we use whole genome resequencing aligned to the Nile tilapia genome assembly

95  (Brawand et al., 2014; Conte et al., 2019) to identify species-informative SNPs distinguishing

96  native and introduced tilapia species important to aquaculture in Tanzania. We also use the SNP

97  panel to identify hybrids in wild populations showing improvement on previous phenotypic and

98  molecular methods for species identification.

99

## Materials and Methods

### Sample collection

102     Samples of 12 *Oreochromis* species in Tanzania were collected by experimental seine

103  netting or purchasing fish directly from fish markets or landing sites. Voucher specimens were

104  stored in 80% ethanol (with photographs taken to record live coloration before preservation), and

105  fin clips for genetic analysis were preserved in 96-100% ethanol or DMSO salt buffer.

106     Specimen ID and sample collection localities are detailed in Table S1. Specimens were

107  identified to species level in the field based on phenotype following diagnostic criteria (Genner et

108  al., 2018). Putative purebred or hybrid status was estimated by a consensus of experienced field

109  researchers from colour photos taken in the field upon collection. Phenotypically identified

110  hybrids either had intermediate phenotypic traits, or discordant combinations of traits typical of

111  purebreds. Species assignments from the genetic data were also compared against this

112  phenotypic ID.

113

## SNP panel design material

115        A set of 25 reference individuals from four species were used to identify optimal SNPs for

116     the panel. These reference individuals were from putatively pure populations of *O. urolepis*

117     (n=10) from the Lower Wami and Rufiji rivers, *O. niloticus* (n=6) and *O. leucostictus* (n=6) from

118     Lake Albert and Lake Victoria, as well as samples from *O. shiranus* (n=3; two from Lituhu and

119     one from aquarium stock at Bangor University) (Figure 1; Table S1). These individuals were

120     classified as reference based on a lack of hybrids or other species in the sampling location and

121     confident morphological identification. A further 75 individuals from an additional eight species

122     were included for joint genotyping, as well as testing the ability of the SNP panel to distinguish

123     species not involved in panel design. Collectively these 100 individuals are referred to herein as
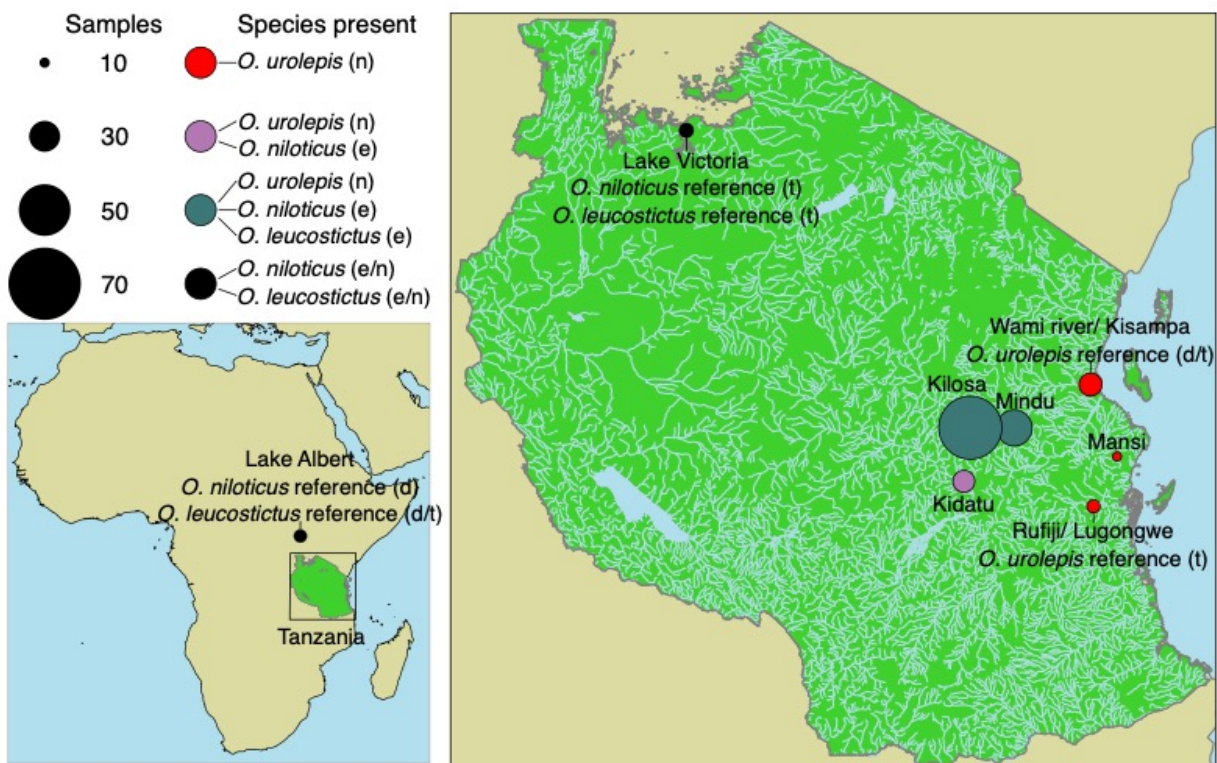
124     the "panel design dataset"



125

126 **Figure 1.** Sample locations for the three focal species within Tanzania (right panel) and Lake
127 Albert (Uganda; bottom left). Abbreviations in 'Species present': n: native; e: exotic.
128 Abbreviations in reference notations on map: d: genome-wide sequencing used for design of
129 SNP array; t: test individuals sequenced using SNP array, used as reference for assigning
130 species. Shapefiles sourced from the ArcGIS Hub (continental boundaries), the ICPAC
131 GeoPortal (Tanzania rivers) and the Humanitarian Data Exchange (Tanzania boundary).
132

## SNP panel performance test material

134     To study the performance of the SNP panel in hybrid identification, we analyzed samples

135 collected during Feb 2015 and May 2016 from i) the Mindu Dam on the Ruvu River near

136 Morogoro (ii) sites near Kilosa on the Wami catchment, iii) the Kidatu reservoir on the Great

137 Ruaha River – a tributary of the Rufiji system, and iv) sites near Utete on the floodplain of the

138 lower Rufiji River, including the oxbow Lake Lugongwe (Figure 1; Table S1). This is

139 subsequently referred to herein as the "panel test dataset". The native species at all four

140 locations is *O. urolepis*, also known in aquaculture literature as *O. urolepis hornorum* or *O.*

141 *hornorum* (Trewavas 1983). Previous work indicates that the Lugongwe site contains pure *O.*

142 *urolepis*, while the introduced *O. niloticus* is established at Kidatu, and both *O. niloticus* and the

143 non-native *O. leucostictus* are present at Mindu and Kilosa (Shechonge et al., 2018).

144 Microsatellite analysis suggested that hybridization was occurring between both introduced

145 species and *O. urolepis* at Mindu and between *O. niloticus* and *O. urolepis* at Kidatu

146 (Shechonge et al., 2018). Reference individuals were included within the panel test dataset (not

147 the same reference individuals as in the panel design dataset). Specifically, six *O. niloticus,* eight

148 *O. leucostictus* and 14 *O. urolepis* individuals were identified based on a lack of hybrids or other

149 species in the sampling location and confident morphological identification.

150

## DNA extraction, whole genome resequencing, read mapping and

## variant calling

7

153    DNA for whole genome resequencing was extracted from fin clips using a PureLink®

154    Genomic DNA extraction kit (Life Technologies). DNA extractions for SNP genotyping were

155    processed using the PureLink Genomic DNA kit or a high-salt extraction protocol. Genomic

156    libraries for paired-end sequencing on the Illumina HiSeq 2500 machine were prepared

157    according to Illumina TruSeq HT protocol to obtain paired-end reads, by the Earlham Institute

158    Genomic Pipelines team. For the 100 individuals in the panel design dataset, low coverage

159    sequencing (target 5X mean depth per sample) were sequenced on the Illumina HiSeq 2500

160    using version 4 chemistry and a 125bp paired-end reads. For the 35 samples in the panel test

161    dataset, sequencing was instead on the Illumina NovaSeq 6000, using 150bp paired-end reads,

162    with an average mean depth of 9x per sample. Raw reads will be made available in the

163    European Nucleotide Archive upon acceptance for publication.

164    For the panel design dataset, quality analysis of raw reads was carried out using fastQC

165    (v0.11.1) (Andrews, 2010). Alignment and duplicate removal were conducted using a local

166    (Earlham Institute) instance of the Galaxy platform (Blankenberg et al., 2010; Giardine et al.,

167    2005; Goecks et al., 2010). Low coverage reads were all aligned to the *Oreochromis niloticus*

168    reference genome [consisting of the NCBI Orenil1.1 genome version GCA_000188235.2

169    (Brawand et al., 2014), concatenated with the NCBI mitochondrial genome GU238433.1], using

170    the default settings of BWA-MEM (Galaxy tool version 0.7.12.1) (Li, 2013). Duplicates were

171    removed using the samtools (Galaxy tool version 2.0) rmdup tool (Li et al., 2009). Local

172    realignment around indels was performed per sample using the IndelRealigner tool from

173    software package GATK v 3.5.0 (McKenna et al., 2010). A reference sequence dictionary was

174    created for the reference file using PicardTools v.1.140 (http://broadinstitute.github.io/picard),

175    and the index files for the reference and aligned bam files were created using samtools (v.1.3)

176    faidx and samtools index.

177    SNP and short indel variants were called against the reference genome using GATK v

178    3.5.0 Haplotypecaller, using the options -ERC GVCF to output gvcf format, and –minPruning and

179    –minDanglingBranchLength parameters set to 1, to account for the low levels of coverage in the

180    resequencing dataset. Variants were called using a sequence dataset for 100 individuals

181    including pure *Oreochromis* species and putative hybrids. Variant evaluation was performed

182    using PicardTools Collect Variant Calling Metrics function. Output variants were separated by

183    SNP/indel and nuclear/mitochondrial scaffolds, and thereafter analysed separately. Indel files

184    were used to mask indels and sites within 5 bp of indels in the gvcf files,but were otherwise not

185    included in analysis. Variant filtration was performed using GATKs VariantFiltration tool using the

186    following hard filters: QD (QualbyDepth): < 2.0; FS (FisherStrandBias) > 20.0; SOR

187    (StrandOddsRatio)    >    4.0;    MQ    (RMSMappingQuality)    >    40.0;    MQRS

188    (MappingQualityRankSumTest) < -2.5; RPRS (ReadPosRankSumTest) < -2.0.

189    For the panel test whole-genome resequence data, variants were called against a newer

190    version of the *Oreochromis niloticus* reference genome (Conte et al., 2019; GCF_001858045.2),

191    not including any of the individuals used to design the panel (Table S2). SNPs were called as for

192    the SNP array design whole-genome calls, except for slightly different filtering parameters (sites

193    excluded   with   Quality-by-depth   <   2,   FS   >   60,   MQ   <   40.0,   MQRankSum   <   -12.5,

194    ReadPosRankSum < -8 or total depth less than 100 or greater than 3000). Unlike the SNP array

195    design genotype calls, the galaxy toolkit was not used for this dataset, and different versions of

196    bwa (v0.7.17) and samtools (v1.10) were used. Using bcftools (v1.10.2), biallelic SNPs with a

197    minor-allele count of at least three were extracted and pruned for missing taxa less than 50%

198    and linkage using the prune function, removing SNPs with $R^2$ greater than 0.6 over 50kb

199    windows.

200

## Identification of SNPs for the panel

201

202     Biallelic nuclear SNPs from the panel design SNP set (only aligned to linkage groups and

203     excluding those mapped to unplaced scaffolds) were extracted, and the dataset was filtered to

204     include only the 25 reference individuals. Vcftools (v0.1.13) (Danecek et al., 2011) was used to

205     calculate pairwise $F_{ST}$ values (-vcf-weir-pop) between each of the reference species groups (*O.*

206     *urolepis* n=10; *O. niloticus* n=6; *O. leucostictus* n=6; *O. shiranus* n=3). The SNP set was filtered

207     to include pairwise $F_{ST}$ values >0.9 for at least three of the six pairwise reference population

208     comparisons. The SNP list was further filtered by imposing a minimum distance of 2mn bp

209     between SNPs and ensuring an even spread of high $F_{ST}$ comparisons across all linkage groups

210     (Table S3).

211     To examine how the SNP set performs in resolving species, the SNPs included in the

212     panel were extracted from the vcf file for all 100 individuals from all 12 study species. Principal

213     Components Analysis (PCA) was then carried out using SNPRelate (Zheng et al., 2012) in R (R

214     Core Team 2019) and plotted using ggplot2 (Wickham, 2016). A neighbor-joining tree was also

215     inferred using the 'ape' package in R (Paradis et al., 2004), using genetic distances calculated

216     using VCF2Dis (https://github.com/BGI-shenzhen/VCF2Dis; accessed December 2020).

217     As the initial analysis and SNP panel design was conducted using an older version of the

218     *O. niloticus* genome assembly, coordinates were subsequently converted to the latest version of

219     the *O. niloticus* reference genome (Conte et al., 2019; GCF_001858045.2) using the NCBI

220     remap tool (https://www.ncbi.nlm.nih.gov/genome/tools/remap: accessed August 2020).

221     Coordinates for both versions of the reference genome are given in Table S4.

222

## SNP panel sequencing

223

224    The selected SNPs were prepared for panel design by extracting a 50bp flanking

225    sequence either side of the SNP locus from the reference genome assembly. Agena

226    Bioscience® (San Diego, California) SNP genotyping of the selected 120 SNPs was performed

227    at the Wellcome Sanger Institute for n=164 samples (see Table S1). This included the remaining

228    reference individuals not used for SNP panel design as well as all the test individuals. Primer

229    and probe sequences for the genotyping are given in Table S4.

230

## SNP panel downsampling

232    To test how many SNPs from the panel are necessary to accurately detect hybrids and

233    assign species, we generated 100 replicates of random subsets of 10, 20, 30, 40, 50, 60, 70, 80,

234    90, 96, 100 and 110 SNPs. We also tested an optimum set of 96 SNPs, according to principal

235    component loading scores, which we calculated for each SNP using PLINK (v1.90) (Purcell et

236    al., 2007). We selected the 96 SNPs with the highest absolute loading values for PC1 and PC2

237    combined. We tested the log-likelihood of fastSTRUCTURE (v1.0.0) runs from $K$=1-12 for each

238    replicate, the number of hybrids (individuals with no ancestry component > 80% identified by

239    fastSTRUCTURE, following Shechonge et al., (2018), for each replicate, and the variability

240    between replicates with the same number of SNPs in ancestry component of each identified

241    hybrid. The optimal 96 SNP set was also compared with the full SNP set using fastSTRUCTURE

242    from $K$=1-12.

243

## SNP panel comparison to other datasets

245    As genotyping is frequently performed in 96-well plates, the optimum panel of 96 SNPs

246    (described above) was compared to the full 120 SNPs to compare performance. Results of the

247    96 SNP panel genotyping were also compared to whole genome resequencing analysis, existing

248    published microsatellite data (Shechonge et al., 2018) and the morphological identification.

249    Microsatellite data was available for 54 of the same individuals used for the SNP panel, and

250    whole genome data was available for 35 of the same individuals used for the SNP panel (see

251    Table S1). No individuals had data available for all three comparisons. The full genome test

252    individual dataset was as described earlier.

253

254    **Population structure and hybrid detection**

255        For the 96 SNP panel and full genome test individual datasets, we performed a Bayesian

256    clustering analysis in the program fastSTRUCTURE (v1.0), running the main algorithm with $K$=1-

257    12. The optimal $K$ value was chosen using the ChooseK.py script within fastSTRUCTURE.

258        For the microsatellite dataset, STRUCTURE (v2.3.4) (Pritchard et al., 2000) was run with

259    500,000 iterations, following 250,000 burn-in iterations. Following (Shechonge et al., 2018), prior

260    cluster assignments (using *LOCPRIOR*) were used, identified using the find.clusters algorithm

261    within the R package adegenet (Jombart, 2008), retaining 20 PCA axes and using 1000

262    iterations. Three clusters ($K$=3) were utilized, corresponding to the three sampled species. Ten

263    independent runs carried out at $K$=3, with the run with the lowest log likelihood utilized to

264    compare to other datasets. The find.clusters algorithm of adegenet was separately run without

265    specifying the number of clusters, to check if the optimal number of clusters according to BIC

266    score differed from the number of species. Additionally, the analysis of microsatellite data was

267    repeated without prior assignments, with 10 replicates for each value of $K$=1-7. The web version

268    of STRUCTURE HARVESTER (Earl & vonHoldt, 2012) was used to infer the most likely value of

269    $K$ using $\Delta K$ (Evanno et al., 2005).

270        Hierarchical clustering results are influenced by the numbers of individuals sampled in

271    each population (Puechmaille, 2016). To prevent this being a confounding factor for the 96 SNP

272   set versus microsatellite and 96 SNP set versus genome-wide comparisons, the 96 SNP set was

273   pruned to only include the relevant individuals for both comparisons. On each of these

274   subsampled datasets, three independent runs of fastSTRUCTURE were carried out for each

275   value of *K* between 1 and 12. Therefore three separate fastSTRUCTURE analyses were

276   performed for the 96 SNP set; one comprising all individuals, one with only the same individuals

277   as in the microsatellite dataset and another with only the same individuals as the genome-wide

278   dataset.

279        For each of these analyses, results from one run with the optimal *K* value were used to

280   assign individuals to species. For the microsatellite analysis with *LOCPRIOR* assignments, the

281   run of *K*=3 with the best log-likelihood score was used. The reference individuals were used to

282   classify ancestry components to species. For the test specimens, a threshold of 80% ancestry

283   component was used to designate individuals to a species (Shechonge et al., 2018), and were

284   considered to have a significant ancestry component corresponding to a species if they had at

285   least 20% ancestry component corresponding to it. For example, if an individual had an ancestry

286   component of 67% corresponding the ancestry component found in the reference *O.*

287   *leucostictus* individuals, and an ancestry component of 33% corresponding to the reference *O.*

288   *urolepis* individuals, it was designated as a *O. leucostictus* x *O. urolepis* hybrid. If it had an

289   ancestry component of 81% corresponding to the *O. leucostictus* reference individuals, and 19%

290   to the *O. urolepis* individuals, it was classified as a *O. leucostictus*.

291        To further assess the hybrid status of individuals described as hybrid from the

292   fastSTRUCTURE analysis (two ancestry components > 0.2), we used NewHybrids (v1.1)

293   (Anderson & Thompson, 2002). NewHybrids assesses the posterior probability that an individual

294   comes from one of six classes: nonhybrid of either of two populations, F1 hybrid, F2 hybrid or

295   backcross of either of two populations. As NewHybrids assumes that there are only two parental

296   taxa, we analysed separate datasets consisting only of individuals assigned by

297 fastSTRUCTURE to belong to one of two species, or to be a hybrid between the two species.

298 For each two-species comparison, five independent runs were carried out, each with a burn-in

299 length of 50,000 followed by an MCMC length of 100,000. No prior information was used to

300 designate individuals to either population. We then checked whether all runs converged (less

301 than 0.05 difference between maximum and minimum estimates in posterior probability for each

302 category for each individual between the five runs) and took the mean posterior probability for

303 each.

304

305 # Results

306 **SNP panel design**

307 For each of the 100 specimens used to generate the SNP panel, ≥98% of reads aligned

308 to the reference genome, with ≥80% of reads properly paired (Table S2). Following the GATK

309 pipeline and filtering, 29,657,078 biallelic SNPs were called. This was pruned down to

310 18,590,392 SNPs with only the 25 reference specimens, including only sites located within the

311 22 linkage groups with missing data from at most one individual. A set of 4,789 SNPs with a

312 pairwise Fst of least 0.9 in three of the six pairwise reference population comparisons was

313 extracted from this set. The 120 SNP set was extracted from these 4,789 SNPs after pruning by

314 distance. These 120 SNPs had an average pairwise Fst of 0.47 across the six pairwise

315 comparisons. All except three of the SNPs had an Fst of 1.0 in at least one pairwise comparison

316 (Table S3). These SNPs were distributed across 22 linkage groups and separated by at least 1

317 Mbp (Table S4). Two of these SNPs were subsequently discarded as they failed Agena

318 Biosciences QC during assay genotyping, resulting in an array of 118 SNPs.

319 PCA of the 118 SNP set extracted from the vcf file with all 100 individuals suggested that

320 species could be distinguished, with *O. niloticus*, *O. leucostictus*, and *O. urolepis* particularly

14

321    distinct. Purebred representatives of species used in the design of the panel were distinct, but

322    clustered relatively tightly in the space within the first two PC axes. They also largely formed

323    monophyletic clusters in a neighbor-joining tree inferred from the 118 SNPs, with the exclusion

324    of potential hybrid individuals, for example the sample T3D7, which was morphologically

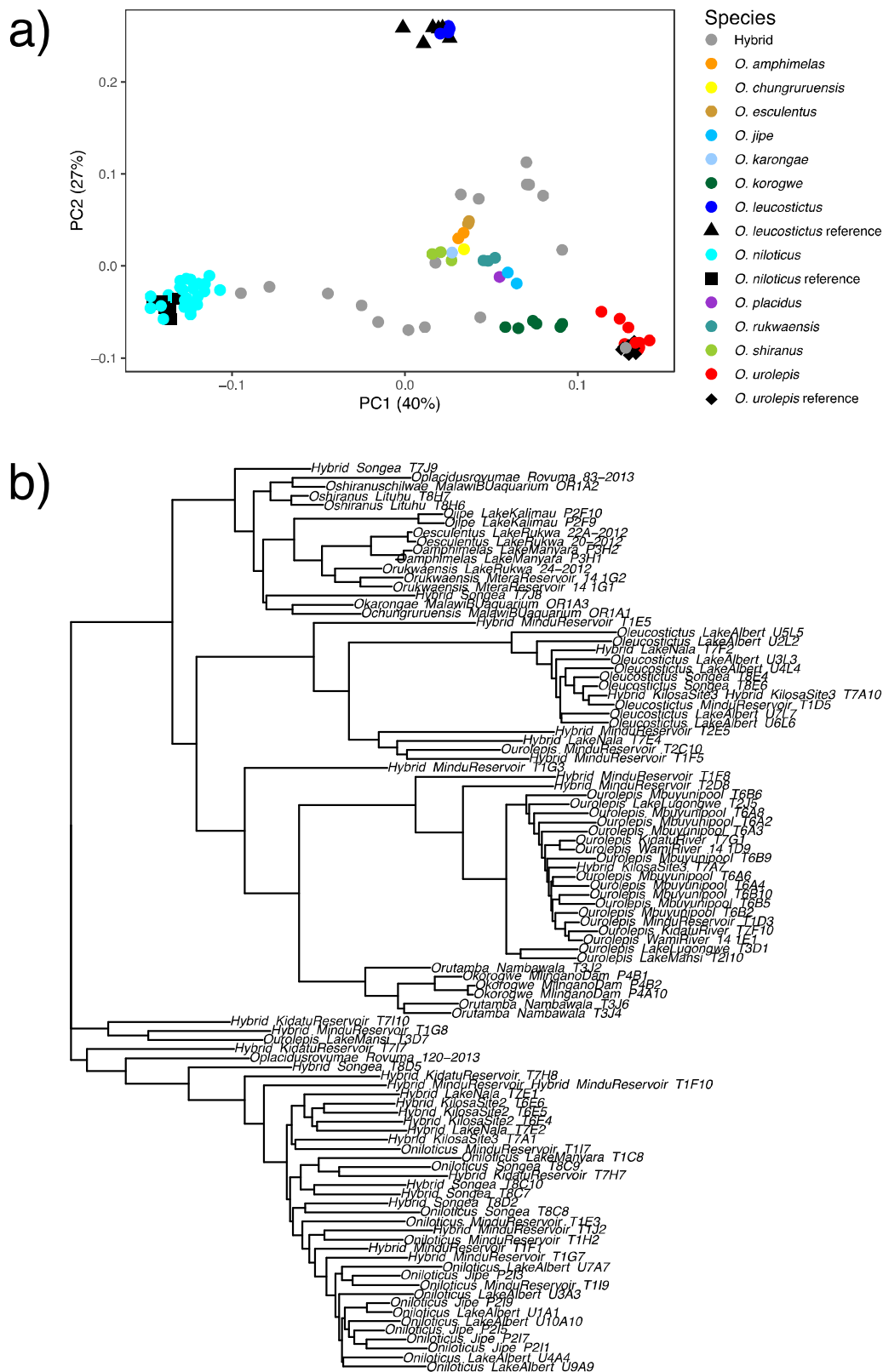325    identified as *O. urolepis* but clustered with morphologically identified hybrids (Figure 2b).

326

327  **Figure 2.** a) PCA of the 118 SNPs, extracted from the full-genome SNP calls from the 100
328  individuals which were used for initial SNP calling to design the SNP panel. b) neighbor-joining
329  tree of the 118 SNPs from the same 100 individuals. Samples are labelled with their
330  morphological ID, sampling location and sample ID, separated by double underscores.
331

## 332  The 96 SNP panel is consistent with full-genome data

333  The 96 SNP set (see Figure S1 for runs from $K$=2-5) and 118 SNP set gave identical

334  classifications for all 164 individuals tested at $K$=3 (the optimal $K$ value for both according to

335  chooseK.py) using fastSTRUCTURE (Figure 3a; Table 1).

336  **Table 1.** Comparison between species assignments between the 96 SNP panel and other
337  datasets.

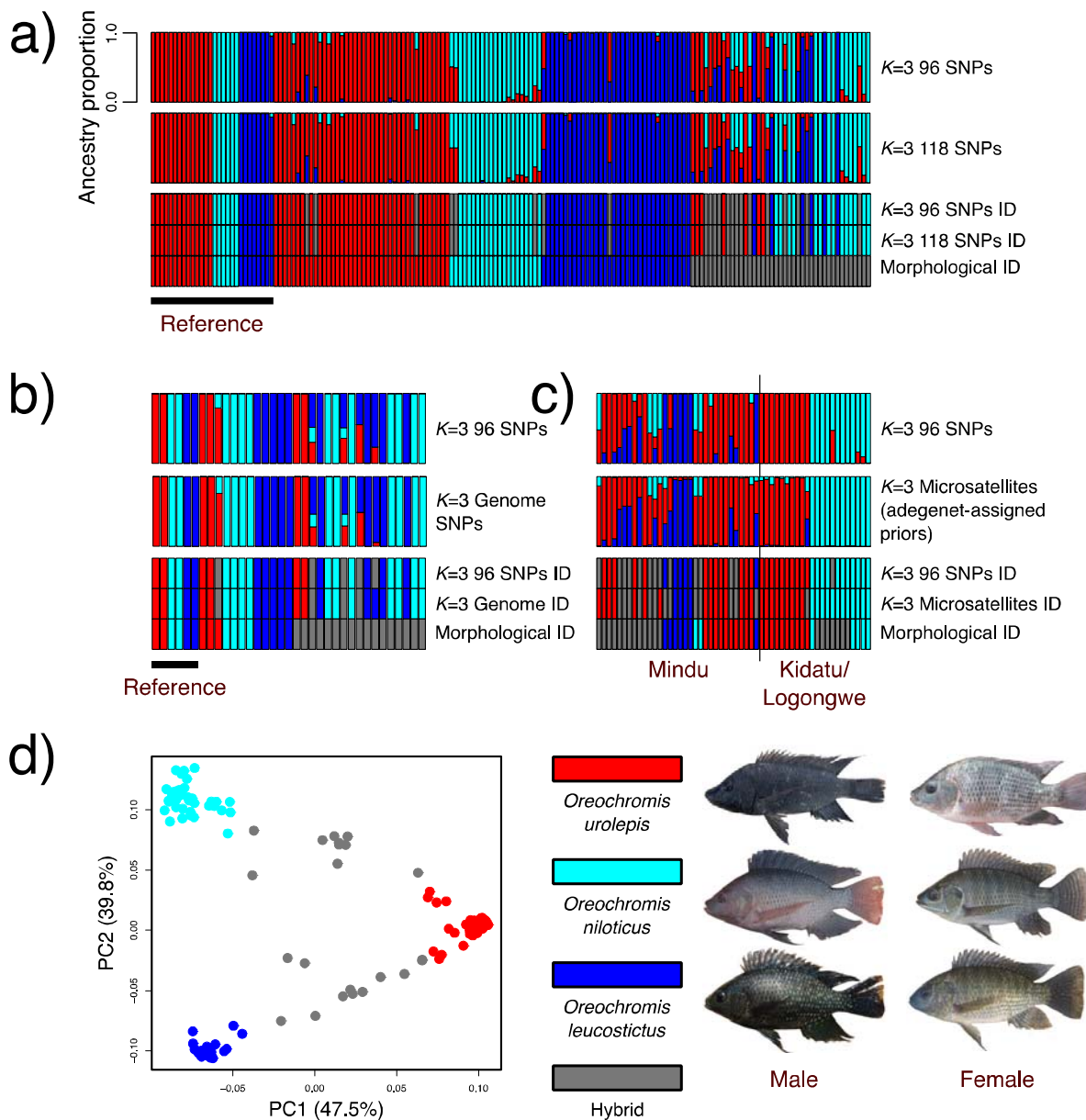| Dataset 1 | Dataset 2 | Number of individuals | Number of individuals with the same assignment | Figure |
|---|---|---|---|---|
| 96 SNP panel | 118 SNP panel | 164 | 164 (100%) | 3a |
| 96 SNP panel | 18 microsatellite array - prior cluster assignment using adegenet | 54 | 48 (89%) | 3c |
| 96 SNP panel | 1,822,719 genome-wide SNPs | 35 | 34 (97%) | 3b |
| 96 SNP panel | Morphological classification | 164 | 129 (79%) | 3a |
| Morphological classification | 18 microsatellite array - prior cluster assignment using adegenet | 54 | 32 (59%) | 3c |
| Morphological classification | 1,822,719 genome-wide SNPs | 35 | 20 (57%) | 3b |

338

17

**Figure 3.** a-c) fastSTRUCTURE analysis, comparison between the 96 SNP set and: a) 118 SNPs; b) genome-wide SNPs; c) microsatellites. d) PCA of the optimal 96 SNP panel, PC1 vs. PC2. The right-hand panel includes representative photographs of mature adults of each species (not to scale).

Following pruning for linkage, 1,822,719 SNPs were used for the genome-wide fastSTRUCTURE analysis. For each of five independent runs $K=3$ was identified as the value to

18

347 both maximize marginal likelihood and explain the structure in data. Out of the 35 individuals

348 with both genome-wide and 96 SNP data, 34 were identified consistently between the two

349 (Figure 3c; Table 1). The only individual they differed on was identified as *O. leucostictus* in the

350 genome-wide data, whereas the 96 SNP set identified it as a hybrid, with a majority *O.*

351 *leucostictus* ancestry but also a *O. urolepis* component. Phenotypically, 17 of these individuals

352 were identified as hybrid. However, only 3 were identified as hybrid in the genome-wide data

353 and 4 in the 96 SNP set. Seven of the phenotypic hybrids were identified as *O. niloticus*, two as

354 *O. urolepis* and five as *O. leucostictus* in the 96 SNP set, with four on the genome-wide data.

355 **The 96 SNP panel outperforms microsatellites**

356 In the STRUCTURE analysis of 18 microsatellites, using the method of Shechonge et al.,

357 (2018), where prior assignment of specimens to clusters based on the known number of species

358 ($K$=3), 89% of individuals were given the same assignment as the 96 SNP panel (Table 1;

359 Figure 3c).

360 However, adegenet incorrectly identified $K$=4 as the optimal number of clusters,

361 according to BIC score. Equally, STRUCTURE analyses with no *a priori* clustering of specimens

362 suggested an optimal $K$ value of $K$=2, with another small peak at $K$=5 (Figure S2). In general,

363 assignments without *a priori* information gave an unclear pattern, distinguishing *O. niloticus* and

364 *O. urolepis* from Kidatu and Lugongwe,,but failing to distinguish species in Mindu (Figure S3).

365 The 96 SNP panel, based on the same samples, by contrast, consistently identified $K$=3 as the

366 optimal number of clusters, and reliably distinguished species in Mindu (Figure S3).

367 **The 96 SNP panel is more accurate than identification based on**

368 **phenotype**

19

369    The 96 SNP set assigned all the reference individuals to the same species as the

370    phenotypic ID (Table S5). However, there were some differences in the test individuals of all

371    three species, with three phenotypically identified *O. urolepis,* three phenotypically identified *O.*

372    *niloticus* and two phenotypically identified *O. leucostictus* being designated as hybrids (all *O.*

373    *niloticus* x *O. urolepis* or *O. leucostictus* x *O. urolepis*). One phenotypically identified *O.*

374    *leucostictus* was instead classified as *O. urolepis* by the 96 SNP panel. Many of the

375    phenotypically identified hybrids were instead given pure species classification: six as *O.*

376    *urolepis*, 14 as *O. niloticus* and six as *O. leucostictus*. Only 14 were identified as hybrid by both

377    phenotype and the 96 SNP panel (Figure 3a).

378

## Validation of hybrid classification

380    In total, 96% of fastStructure identifications were corroborated with NewHybrids, with

381    posterior probability > 0.98 (Table S6,7). Two NewHybrids analyses were carried out: one with

382    *O. urolepis* and *O. niloticus* individuals, and individuals identified as hybrid between the two, and

383    one with *O. urolepis* and *O. leucostictus* individuals, and their hybrids. Comparisons were not

384    made between *O. niloticus* and *O. leucostictus,* as no hybrids were identified between the two

385    using fastSTRUCTURE. Six F1 hybrids were identified between *O. urolepis* and *O. niloticus*

386    (posterior probability > 0.95), alongside five *O. urolepis* backcrosses and three *O. niloticus*

387    backcrosses (Table S6). Five F1 hybrids (posterior probability > 0.95) were identified between

388    *O. leucostictus* and *O. urolepis*, alongside one F2 hybrid. Six *O. urolepis* backcrosses were also

389    identified, with three *O. leucostictus* backcrosses (Table S6).

390    Together, we found evidence of introgression between the native *O. urolepis* and both

391    the invasive *O. leucostictus* and *O. niloticus* in Kilosa, Kidatu and Mindu. We find no evidence of

392    hybrids in Rufiji, Lugongwe or Mansi (Table S1). We found no evidence of any hybrids between

393    *O. niloticus* and *O. leucostictus*. See Table 2 for the numbers of each species and hybrid

394    identified by the 96 SNP panel.

395    **Table 2**. Number of individuals of each species identified in each sampling location by the 96
396    SNP panel.

| Location | *O. urolepis* individuals | *O. leucostictus* individuals | *O. niloticus* individuals | Hybrids |
|---|---|---|---|---|
| Mindu Dam | 13 | 6 | 0 | 7 *leucostictus* x *urolepis*<br>6 *niloticus* x *urolepis*<br>1 *leucostictus* x *niloticus* x *urolepis* |
| Kilosa | 6 | 32 | 18 | 4 *leucostictus* x *urolepis*<br>1 *niloticus* x *urolepis*<br>1 *leucostictus* x *niloticus* x *urolepis* |
| Kidatu resevoir | 5 | 0 | 14 | 2 *niloticus* x *urolepis* |
| Lake Mansi/ Lugongwe/ Mindu | 20 | 0 | 0 | 0 |

397

398    **Different subsets of at least 80 out of the 118 SNP dataset give**

399    **consistent results**

400    For the majority of subsample replicates of the 118 SNPs, *K=3* was identified as the

401    model complexity that maximized marginal likelihood for the majority of these replicates with the

402    following exceptions: *K=2* was optimal for 6/100 10-SNP sets, *K=4* was chosen for 1/100 20-

403    SNP sets, 1/100 of the 30-SNP sets and 1/100 of the 70-SNP sets, and *K=5* was chosen for

404    1/100 of the 60-SNP sets and 1/100 of the 90-SNP sets. The Model components used to explain

405    structure in the data varied much more between replicates, from 1-11.

406        Increasing the number of SNPs increased likelihoods, with sharp increases from 10-30

407    SNPs and more modest increases thereafter (Figure S4a). The number of iterations in which all

408    the reference individuals were correctly classified into populations increased with the number of

409    SNPs, until it reached 100% at 80 SNPs (Figure S4b). It also decreased the number of hybrids

410    identified up to 80 SNPs, after which it then stabilized (Figure S5a). Increasing SNP number also

411    increased the stability of the estimated hybrid ancestry proportion, measured as the variability in

412    the minor ancestry component of a hybrid (Figure S5b). In the 96 SNP iterations, 18 individuals

413    were consistently identified as hybrid in all replicates, whereas 9 were sometimes classified as

414    hybrids. Of these, 5 were identified in fewer than 12 out of 100 replicates, and 4 were identified

415    in at least 78 replicates (Figure S5c,d).

416

# 417    Discussion

418    We demonstrate that a reduced panel of 96 genome-wide SNPs performs comparatively well to

419    full genome resequencing in distinguishing species and identifying hybrids of *Oreochromis*. We

420    identify replicate cases where introduced aquaculture species have become established and

421    interbred with native species, including backcrosses as well as F1 and F2 hybrids. We

422    demonstrate that hybridization is persistent in the environment with multi-generation hybrids and

423    backcrossing to parental species.

424        We found that the ability of a reduced 96-SNP panel to detect hybrids was

425    indistinguishable from the full 118 SNPs that we genotyped. This is likely to be the most cost-

426    efficient panel size, as genotyping is frequently performed in 96-well plates. None of these SNPs

427    overlapped with previously identified species-diagnostic SNPs for *O. niloticus* (Syaifudin et al.,

428 2019), likely because our SNP set was optimized by interspecific rather than intraspecific

429 variation. Our analyses indicate that the reduced 96 SNP panel can accurately identify the

430 hybrids between *Oreochromis* species tested, including introgression from the invasive *O.*

431 *niloticus* and *O. leucostictus* into the native *O. urolepis* (Shechonge, Ngatunga, Bradbeer, et al.,

432 2019). No hybrids were identified between *O. leucostictus* and *O. niloticus* in Tanzania. These

433 two species co-occur in Lake Albert, Uganda, where they are not known to hybridize (Trewavas

434 1983). It is possible therefore that behavioral, ecological or genomic incompatibilities prevent the

435 two species from hybridizing in populations where they naturally occur, although *O. leucostictus*

436 has been shown to hybridize in Kenya with other subspecies of *O. niloticus,* with which it does

437 not naturally co-exist (Ndiwa et al., 2014).

438 This detection of introgression between *O. urolepis* and *O. leucostictus*, and between *O.*

439 *urolepis* and *O. niloticus* was concordant with previous studies using microsatellite data

440 (Shechonge et al., 2018). Our re-analysis with the same set of microsatellites only gave

441 comparable results if '*LOCPRIOR*' assignments were used based on an initial clustering, which

442 suggests that power was low in the microsatellite analysis due to a small number of markers or

443 samples (Porras-Hurtado et al., 2013). Additionally, using the '*LOCPRIOR*' required choosing

444 the value of *K* based on sampling (*K*=3), rather than the optimal number according to BIC score

445 of *K*=4. This suggests that an added benefit of the 96 SNP set is that prior assumptions are not

446 necessary to set the appropriate value of *K* when relatively few individuals are sampled, unlike

447 with the microsatellite data. This may be important in cases where an unknown number of test

448 species are sampled, or there is hidden population structure (Porras-Hurtado et al., 2013). The

449 SNP panel would therefore require less thorough sampling to allow accurate species or hybrid

450 assignment.

451 Notably, our analyses suggested that morphological identification of hybrids was

452 inconsistent with genetic assignments; many individuals phenotypically assigned as hybrids

23

453    were genetically classified as pure species. This may reflect high phenotypic diversity within

454    species (Table S5), and possibly overlap in characteristics between species, which could be

455    difficult to catalogue, making species identification more challenging. It may also reflect

456    introgression which has been masked by several generations of backcrossing. This would result

457    in small ancestry components for the introgressed species, and incorrect pure species

458    assignment using hierarchical clustering (e.g. STRUCTURE or fastSTRUCTURE) or

459    NewHybrids. Further studies with large sample sizes, thorough population sampling, genomic

460    data and detailed demographic analyses are necessary to identify if this is the case. This would

461    mean introgression has been occurring for several generations, possibly influencing phenotypic

462    variation within species.

463    Several analyses suggested that the panel of 96 SNPs provides sufficient power to

464    reliably identify these species and hybrids. Importantly, species assignments using the 96 SNP

465    panel were almost identical to those given by genome-wide data (Figure 3b; Table 1). This

466    suggests that adding more SNPs at extra cost would not greatly improve assignment accuracy,

467    and introgression from the invasive species can be detected reliably without the considerable

468    investment of whole-genome or reduced-representation (e.g. RAD-seq) resequencing. SNP

469    panels of similar sizes have proved accurate at detecting hybrid status and introgression

470    between domestic cats and European wildcats (Oliveira et al., 2015), and between farmed and

471    wild Atlantic salmon (Wringe et al., 2019).

472    Subsamples of the full 118 SNP set further indicated that individuals could be accurately

473    assigned if > 80 SNPs were used. Although hybrids identification was not fully consistent

474    between iterations even when a larger number of SNPs were used, variability in ancestry

475    components between iterations was low (<0.1). This indicates that the individuals which are

476    classified as hybrids in some but not all iterations (Figure S5d) are those with an ancestry

477    component of close to our arbitrary cut-off to define a hybrid. We recommend that any

478     individuals which are close to the cut-off value chosen are further investigated, for example

479     using NewHybrids. The choice of threshold to define a hybrid may also be adjusted depending

480     on the application. For example, if the panel is being used to eliminate hybrids from breeding

481     stock, then it may be necessary to use a stricter threshold to define hybrids. These analyses

482     indicate that 96 SNPs is above the point of diminishing returns for hybrid identification and

483     accurate reference individual identification, meaning that even if some SNPs fail to amplify in

484     some individuals there should still be sufficient power.

485         It is important to further consider methodological limitations to accurate species and

486     hybrid assignment using the SNP panel. A signal of introgression indicated by hierarchical

487     clustering can be given in the absence of any introgression of one population that has

488     undergone a recent bottleneck, or in the case of 'ghost' introgression from an unsampled

489     population (Lawson et al., 2018). Given that introgression was only inferred in some individuals

490     within each population in our analysis, and the general concordance with NewHybrids analysis,

491     it is likely that the signal we are detecting is in fact introgression, rather than a population-level

492     bottleneck. However, this must be a consideration for users applying the SNP panel on other

493     *Oreochromis* species that we have not tested here. The issue of introgression from unsampled

494     taxa is more likely to be confounding in our dataset, given we have only extensively tested three

495     out of the at least 37 species of *Oreochromis* (Ford et al., 2019). As *O. niloticus* and *O.*

496     *leucostictus* are the only introduced *Oreochromis* species found in the tested water bodies

497     (Shechonge, Ngatunga, Bradbeer, et al., 2019), it is likely that these results do reflect

498     introgression from one of these species. Reassuringly, PCA and a neighbor-joining tree inferred

499     from the 118 SNPs extracted from the individuals with full-genome resequencing suggested that

500     most of the species are distinguishable, particularly the highly invasive *O. niloticus* and *O.*

501     *leucostictus* (Figure 2), suggesting that introgression from either of these two species would be

502     identifiable. The discriminatory ability of the SNP panel will need to be tested in cases where

503     other Tanzanian native species co-occur with the focal introduced species, as the current SNP

504     panel was not optimized for other species groups. However, even if native species could not

505     convincingly be distinguished, the SNP panel we present will be able to identify introgression

506     from invasive *O. niloticus* and *O. leucostictus*.

507         Hierarchical clustering results may also be influenced by uneven sampling of populations

508     (Puechmaille, 2016). In the case that only one or two individuals are sequenced from one

509     population in a large dataset, it is unlikely that they will be assigned a distinct cluster, even in the

510     absence of any introgression. This may mean a lot of diversity within the dataset may be missed.

511     Species assignment tools based on network estimation have the potential to identify these

512     'outlier individuals', which do not belong to any of the reference populations (Kuismin et al.,

513     2020). However, it is not clear how they perform in the presence of hybrid individuals. Future

514     studies using this SNP panel will need to prioritize establishing a reference set of individuals

515     belonging to each target species, with a similar number of individuals of each.

516         We anticipate that our efficient SNP panel will be of use to the aquaculture and

517     conservation genetics communities in assessing broodstock purity, determining hybrid status of

518     wild populations, and identifying populations most in need of conservation resources.

519

# 520   Acknowledgements

536 # Data Accessibility

537 Whole genome resequencing data will be deposited before publication

538 SNP datasets: Dryad

539

540

541 # Author Contributions

542 GFT, MJG, FDP, and MM conceived the study. MJG, GFT, AS, BPN, and RT designed fieldwork

543 and sampling. AGPF, NK, BPN, AS, GFT, RT and MJG conducted or supervised fieldwork, or

544 collected data. AGPF and TM performed laboratory work. AGC, AGPF, GE, LP-D and WH

545 designed and performed the analysis. AGC and AGPF wrote the first draft of the manuscript. All

546 authors commented on and edited the final manuscript.

547

548 **Supplementary Tables**

557  **Supplementary Figures**



558
559  Figure S1. fastSTRUCTURE analysis for all individuals in the 96 SNP panel dataset, from $K$=2 to
560  $K$=5.
561

28

562
563
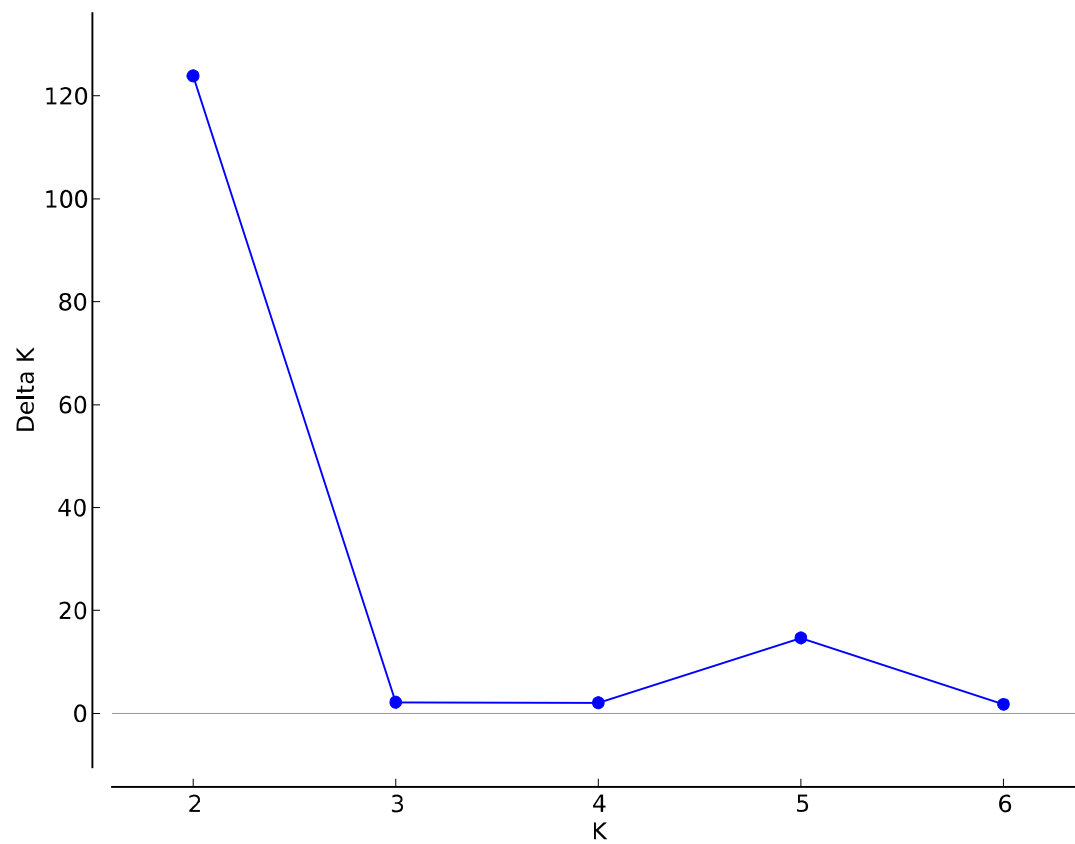564  Figure S2. $\Delta K$ (DeltaK) values for STRUCTURE runs on the microsatellite dataset, without prior
565  assignment, from $K$=2 to $K$=6.
566

567
568
569    Figure S3. Comparison between species assignment using the 96 SNP panel (top row), and
570    microsatellite analysis without prior assignment at *K*=3 & *K*=4.
571



572
573
574    Figure S4. a) Average log-likelihoods for the 100 replicates of each number of sub-sample
575    SNPs. Error bars represent standard error. b) The percentage of replicates for each number of
576    sub-sample where all of the reference individuals were correctly assigned to their species.
577

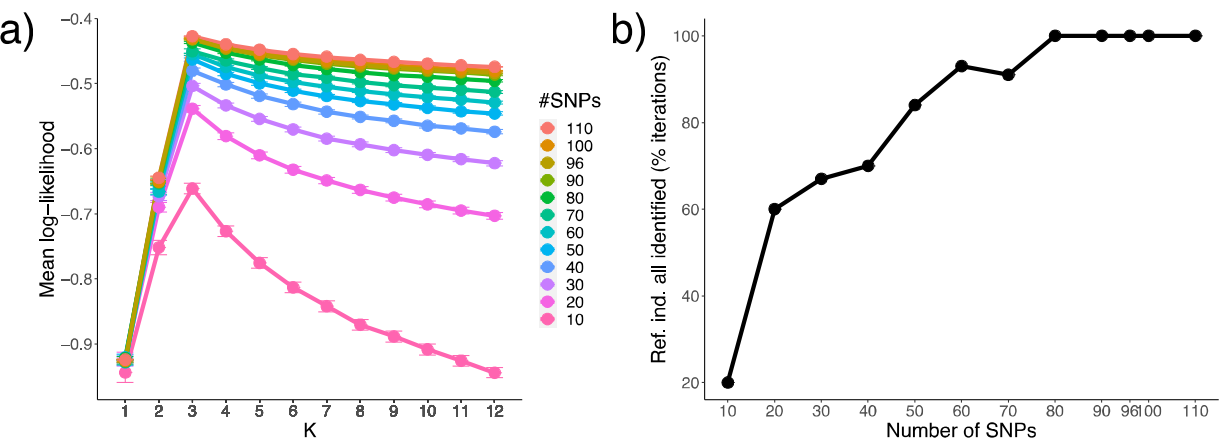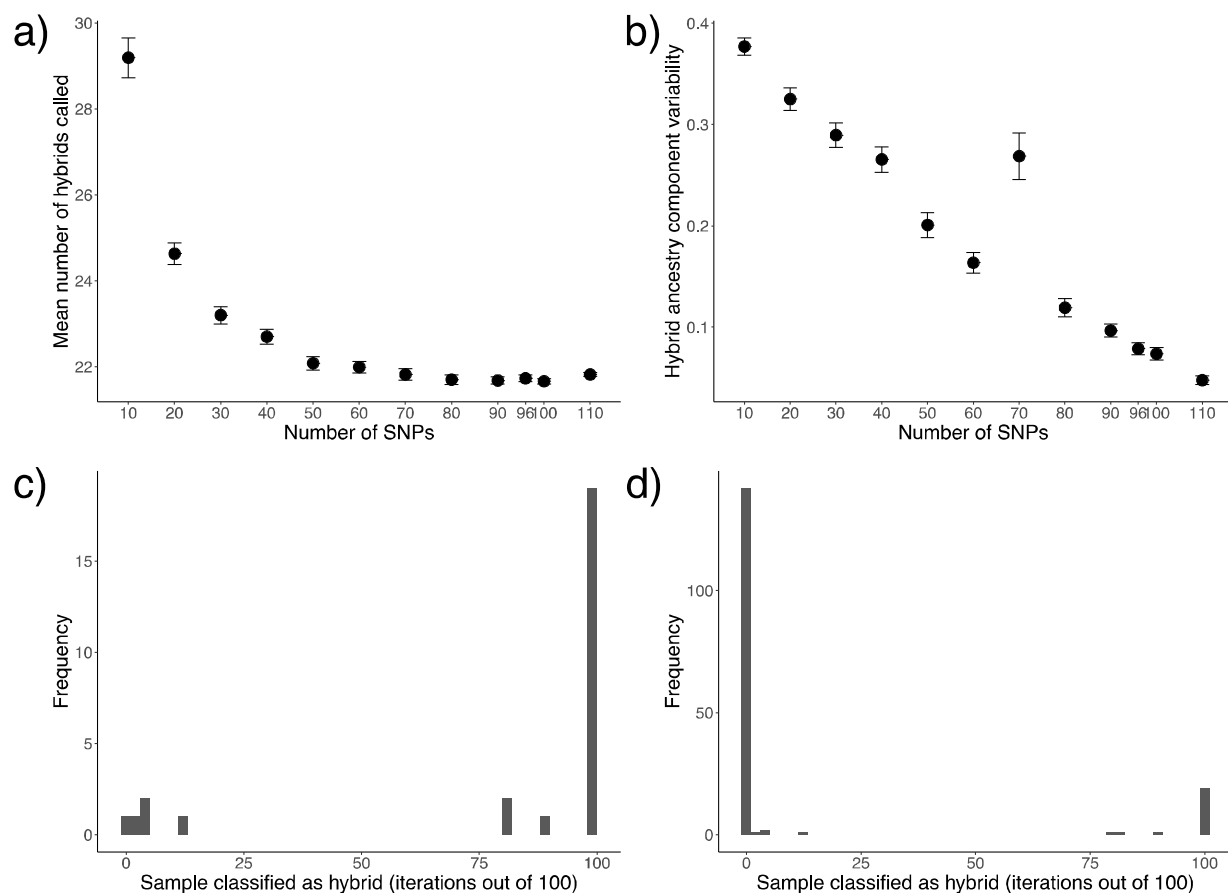Figure S5. a) The mean number of hybrids called (no ancestry component > 80%) across the 100 replicates of each number of sub-sample SNPs. Error bars represent standard error. b) The mean variability in minor ancestry component between the 100 replicates for each individual identified as hybrid in at least one of these replicates for each number of sub-sample SNPs. Error bars represent standard error. c) Histogram of the frequency at which hybrids were classified as hybrids across the 100 replicates of 96 random SNPs. d) same as c) except that data for individuals never identified as hybrids is added.

# References

Anderson, E. C., & Thompson, E. A. (2002). A model-based method for identifying species hybrids using multilocus genetic data. *Genetics*, *160*(3), 1217–1229.
Andrews, S. (2010). FastQC: a quality control tool for high throughput sequence data. https://www.bioinformatics.babraham.ac.uk/projects/ fastqc/.
Angienda, P. O., Lee, H. J., Elmer, K. R., Abila, R., Waindi, E. N., & Meyer, A. (2011). Genetic structure and gene flow in an endangered native tilapia fish (*Oreochromis esculentus*) compared to invasive Nile tilapia (*Oreochromis niloticus*) in Yala swamp, East Africa.

31

596        *Conservation Genetics* , *12*(1), 243–255.

597 Balirwa, J. S. (1988). The evolution of the fishery of *Oreochromis niloticus* (Pisces : Cichlidae) in
598        Lake Victoria. *Hydrobiologia*, *232*(1), 85–89.

599 Bbole, I., Katongo, C., Deines, A. M., Shumba, O., & Lodge, D. M. (2014). Hybridization between
600        non-indigenous *Oreochromis niloticus* and native *Oreochromis* species in the lower Kafue
601        River and its potential impacts on fishery. *Journal of Ecology and The Natural Environment*,
602        *6*(6), 215–225.

603 Blankenberg, D., Kuster, G. V., Coraor, N., Ananda, G., Lazarus, R., Mangan, M., Nekrutenko,
604        A., & Taylor, J. (2010). Galaxy: A Web☐Based Genome Analysis Tool for Experimentalists.
605        *Current Protocols in Molecular Biology / Edited by Frederick M. Ausubel ... [et Al.]*, *89*(1).

606 Bradbeer, S. J., Harrington, J., Watson, H., Warraich, A., Shechonge, A., Smith, A., Tamatamah,
607        R., Ngatunga, B. P., Turner, G. F., & Genner, M. J. (2019). Limited hybridization between
608        introduced and Critically Endangered indigenous tilapia fishes in northern Tanzania.
609        *Hydrobiologia*, *832*(1), 257–268.

610 Brawand, D., Wagner, C. E., Li, Y. I., Malinsky, M., Keller, I., Fan, S., Simakov, O., Ng, A. Y.,
611        Lim, Z. W., Bezault, E., Turner-Maier, J., Johnson, J., Alcazar, R., Noh, H. J., Russell, P.,
612        Aken, B., Alföldi, J., Amemiya, C., … Di Palma, F. (2014). The genomic substrate for
613        adaptive radiation in African cichlid fish. *Nature*, *513*(7518), 375–381.

614 Campbell, N. R., & Narum, S. R. (2011). Development of 54 novel single-nucleotide
615        polymorphism (SNP) assays for sockeye and coho salmon and assessment of available
616        SNPs to differentiate stocks within the Columbia River. *Molecular Ecology Resources*, *11
617        Suppl 1*, 20–30.

618 Canonico, G. C., Arthington, A., McCrary, J. K., & Thieme, M. L. (2005). The effects of
619        introduced tilapias on native biodiversity. *Aquatic Conservation: Marine and Freshwater
620        Ecosystems*, *15*(5), 463–483.

621 Conte, M. A., Joshi, R., Moore, E. C., Nandamuri, S. P., Gammerdinger, W. J., Roberts, R. B.,
622        Carleton, K. L., Lien, S., & Kocher, T. D. (2019). Chromosome-scale assemblies reveal the
623        structural evolution of African cichlid genomes. *GigaScience*, *8*(4).

624 D'Amato, M. E., Esterhuyse, M. M., Van Der Waal, B. C. W., Brink, D., & Volckaert, F. A. M.
625        (2007). Hybridization and phylogeography of the Mozambique tilapia *Oreochromis
626        mossambicus* in southern Africa evidenced by mitochondrial and microsatellite DNA
627        genotyping. *Conservation Genetics* , *8*(2), 475–488.

628 Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., Handsaker, R.
629        E., Lunter, G., Marth, G. T., Sherry, S. T., McVean, G., Durbin, R., & 1000 Genomes Project
630        Analysis Group. (2011). The variant call format and VCFtools. *Bioinformatics* , *27*(15),
631        2156–2158.

632 Deines, A. M., Bbole, I., Katongo, C., Feder, J. L., & Lodge, D. M. (2014). Hybridisation between
633        native *Oreochromis* species and introduced Nile tilapia *O. niloticus* in the Kafue River,
634        Zambia. *African Journal of Aquatic Science*, *39*(1), 23–34.

635 Earl, D. A., & vonHoldt, B. M. (2012). STRUCTURE HARVESTER: a website and program for
636        visualizing STRUCTURE output and implementing the Evanno method. *Conservation
637        Genetics Resources*, *4*(2), 359–361.

638 Evanno, G., Regnaut, S., & Goudet, J. (2005). Detecting the number of clusters of individuals
639        using the software STRUCTURE: a simulation study. *Molecular Ecology*, *14*(8), 2611–2620.

640 FAO. (2020). *The State of World Fisheries and Aquaculture 2020. Sustainability in action*. Rome.

641 Firmat, C., Alibert, P., Losseau, M., Baroiller, J.-F., & Schliewen, U. K. (2013). Successive
642        invasion-mediated interspecific hybridizations and population structure in the endangered
643        cichlid *Oreochromis mossambicus*. *PloS ONE*, *8*(5), e63880.

644 Ford, A. G. P., Bullen, T. R., Pang, L., Genner, M. J., Bills, R., Flouri, T., Ngatunga, B. P., Rüber,

645   L., Schliewen, U. K., Seehausen, O., Shechonge, A., Stiassny, M. L. J., Turner, G. F., &
646        Day, J. J. (2019). Molecular phylogeny of *Oreochromis* (Cichlidae: Oreochromini) reveals
647        mito-nuclear discordance and multiple colonisation of adverse aquatic environments.
648        *Molecular Phylogenetics and Evolution*, *136*, 215–226.
649   Genner, M. J., Connell, E., Shechonge, A., Smith, A., Swanstrom, J., Mzighani, S., Mwijage, A.,
650        Ngatunga, B. P., & Turner, G. F. (2013). Nile tilapia invades the Lake Malawi catchment.
651        *African Journal of Aquatic Science*, *38*(sup1), 85–90.
652   Genner, M. J., Turner, G. F., & Ngatunga, B. P. (2018). A guide to the tilapia fishes of Tanzania.
653        Available at: https://martingenner.weebly.com/tanzania-tilapia-guide.html.
654   Giardine, B., Riemer, C., Hardison, R. C., Burhans, R., Elnitski, L., Shah, P., Zhang, Y.,
655        Blankenberg, D., Albert, I., Taylor, J., Miller, W., Kent, W. J., & Nekrutenko, A. (2005).
656        Galaxy: a platform for interactive large-scale genome analysis. *Genome Research*, *15*(10),
657        1451–1455.
658   Goecks, J., Nekrutenko, A., Taylor, J., & Galaxy Team. (2010). Galaxy: a comprehensive
659        approach for supporting accessible, reproducible, and transparent computational research
660        in the life sciences. *Genome Biology*, *11*(8), R86.
661   Goudswaard, P. C., Witte, F., & Katunzi, E. F. B. (2002). The tilapiine fish stock of Lake Victoria
662        before and after the Nile perch upsurge. *Journal of Fish Biology*, *60*(4), 838–856.
663   Helyar, S. J., Limborg, M. T., Bekkevold, D., Babbucci, M., van Houdt, J., Maes, G. E.,
664        Bargelloni, L., Nielsen, R. O., Taylor, M. I., Ogden, R., Cariani, A., Carvalho, G. R.,
665        FishPopTrace Consortium, & Panitz, F. (2012). SNP discovery using next generation
666        transcriptomic sequencing in Atlantic herring (*Clupea harengus*). *PloS ONE*, *7*(8), e42089.
667   Hess, J. E., Campbell, N. R., Docker, M. F., Baker, C., Jackson, A., Lampman, R., McIlraith, B.,
668        Moser, M. L., Statler, D. P., Young, W. P., Wildbill, A. J., & Narum, S. R. (2015). Use of
669        genotyping by sequencing data to develop a high-throughput and multifunctional SNP panel
670        for conservation applications in Pacific lamprey. *Molecular Ecology Resources*, Vol. 15, pp.
671        187–202.
672   Jombart, T. (2008). *adegenet*: a R package for the multivariate analysis of genetic markers.
673        *Bioinformatics* , *24*(11), 1403–1405.
674   Jorissen, M. W. P., Huyse, T., Pariselle, A., Wamuini Lunkayilakio, S., Muterezi Bukinga, F.,
675        Chocha Manda, A., Kapepula Kasembele, G., Vreven, E. J., Snoeks, J., Decru, E., Artois,
676        T., & Vanhove, M. P. M. (2020). Historical museum collections help detect parasite species
677        jumps after tilapia introductions in the Congo Basin. *Biological Invasions*, Vol. 22, pp. 2825–
678        2844.
679   Kajungiro, R. A., Palaiokostas, C., Pinto, F. A. L., Mmochi, A. J., Mtolera, M., Houston, R. D., &
680        de Koning, D. J. (2019). Population structure and genetic diversity of Nile tilapia
681        (*Oreochromis niloticus*) strains cultured in Tanzania. *Frontiers in Genetics*, *10*, 1269.
682   Kuismin, M., Saatoglu, D., Niskanen, A. K., Jensen, H., & Sillanpää, M. J. (2020). Genetic
683        assignment of individuals to source populations using network estimation tools. *Methods in
684        Ecology and Evolution / British Ecological Society*, *11*(2), 333–344.
685   Larson, W. A., Seeb, J. E., Pascal, C. E., Templin, W. D., & Seeb, L. W. (2014). Single-
686        nucleotide polymorphisms (SNPs) identified through genotyping-by-sequencing improve
687        genetic stock identification of Chinook salmon (*Oncorhynchus tshawytscha*) from western
688        Alaska. *Canadian Journal of Fisheries and Aquatic Sciences. 71*(5), 698–708.
689   Lawson, D. J., van Dorp, L., & Falush, D. (2018). A tutorial on how not to over-interpret
690        STRUCTURE and ADMIXTURE bar plots. *Nature Communications*, *9*(1), 3258.
691   Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
692        *arXiv*, *1307.3997*.
693   Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G.,

694      Durbin, R., & 1000 Genome Project Data Processing Subgroup. (2009). The Sequence
695         Alignment/Map format and SAMtools. *Bioinformatics* , *25*(16), 2078–2079.
696   Lind, C. E., Agyakwah, S. K., Attipoe, F. Y., Nugent, C., Richard P M, & Toguyeni, A. (2019).
697         Genetic diversity of Nile tilapia (*Oreochromis niloticus*) throughout West Africa. *Scientific*
698         *Reports*, 9, 16767.
699   Macaranas, J. M., Taniguchi, N., Pante, M. J. R., Capili, J. B., & Pullin, R. S. V. (1986).
700         Electrophoretic evidence for extensive hybrid gene introgression into commercial
701         *Oreochromis niloticus* (L.) stocks in the Philippines. *Aquaculture Research*, *17*(4), 249–258.
702   Mamonekene, V., & Stiassny, M. L. J. (2012). Fishes of the Du Chaillu Massif, Niari Depression,
703         and Mayombe Massif (Republic of Congo, west-central Africa): A list of species collected in
704         tributaries of the upper Ogowe and middle and upper Kouilou-Niari River basins. *Check List*
705         , *8*(6), 1172–1183.
706   Marufu, L. T., & Chifamba, P. C. (2013). A comparison of diel feeding pattern, ingestion and
707         digestive efficiency of *Oreochromis niloticus* and *Oreochromis macrochir* in Lake Chivero,
708         Zimbabwe. *African Journal of Aquatic Science*, *38*(2), 221–228.
709   McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., Garimella, K.,
710         Altshuler, D., Gabriel, S., Daly, M., & DePristo, M. A. (2010). The Genome Analysis Toolkit:
711         a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome*
712         *Research*, *20*(9), 1297–1303.
713   Mojekwu, T. O., Cunningham, M. J., Bills, R. I., Pretorius, P. C., & Hoareau T. B. (2021). Utility
714         of DNA barcoding in native *Oreochromis* species. *Journal of Fish Biology*, *98*(2), 498-506.
715   Moses, M., Mtolera, M. S. P., Chauka, L. J., Lopes, F. A., de Koning, D. J., Houston, R. D., &
716         Palaiokostas, C. (2020). Characterizing the genetic structure of introduced Nile tilapia
717         (*Oreochromis niloticus*) strains in Tanzania using double digest RAD sequencing.
718         *Aquaculture International: Journal of the European Aquaculture Society*, *28*(2), 477–492.
719   Ndiwa, T. C., Nyingi, D. W., & Agnese, J.-F. (2014). An important natural genetic resource of
720         *Oreochromis niloticus* (Linnaeus, 1758) threatened by aquaculture activities in Loboi
721         drainage, Kenya. *PloS ONE*, *9*(9), e106972.
722   Oliveira, R., Randi, E., Mattucci, F., Kurushima, J. D., Lyons, L. A., & Alves, P. C. (2015).
723         Toward a genome-wide approach for detecting hybrids: informative SNPs to detect
724         introgression between domestic cats and European wildcats (*Felis silvestris*). *Heredity*,
725         *115*(3), 195–205.
726   Paradis, E., Claude, J., & Strimmer, K. (2004). APE: Analyses of Phylogenetics and Evolution in
727         R language. *Bioinformatics* , *20*(2), 289–290.
728   Porras-Hurtado, L., Ruiz, Y., Santos, C., Phillips, C., Carracedo, A., & Lareu, M. V. (2013). An
729         overview of STRUCTURE: applications, parameter settings, and supporting software.
730         *Frontiers in Genetics*, *4*, 98.
731   Pritchard, J. K., Stephens, M., & Donnelly, P. (2000). Inference of population structure using
732         multilocus genotype data. *Genetics*, *155*(2), 945–959.
733   Puechmaille, S. J. (2016). The program structure does not reliably recover the correct population
734         structure when sampling is uneven: subsampling and new estimators alleviate the problem.
735         *Molecular Ecology Resources*, *16*(3), 608–627.
736   Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M. A. R., Bender, D., Maller, J.,
737         Sklar, P., de Bakker, P. I. W., Daly, M. J., & Sham, P. C. (2007). PLINK: a tool set for whole-
738         genome association and population-based linkage analyses. *American Journal of Human*
739         *Genetics*, *81*(3), 559–575.
740   Shechonge, A., Ngatunga, B. P., Bradbeer, S. J., Day, J. J., Freer, J. J., Ford, A. G. P., Kihedu,
741         J., Richmond, T., Mzighani, S., Smith, A. M., Sweke, E. A., Tamatamah, R., Tyers, A. M.,
742         Turner, G. F., & Genner, M. J. (2019). Widespread colonisation of Tanzanian catchments by

743    introduced *Oreochromis* tilapia fishes: the legacy from decades of deliberate introduction.
744    *Hydrobiologia*, *832*(1), 235–253.
745 Shechonge, A., Ngatunga, B. P., Tamatamah, R., Bradbeer, S. J., Harrington, J., Ford, A. G. P.,
746    Turner, G. F., & Genner, M. J. (2018). Losing cichlid fish biodiversity: genetic and
747    morphological homogenization of tilapia following colonization by introduced species.
748    *Conservation Genetics* , *19*(5), 1199–1209.
749 Shechonge, A., Ngatunga, B. P., Tamatamah, R., Bradbeer, S. J., Sweke, E., Smith, A., Turner,
750    G. F., & Genner, M. J. (2019). Population genetic evidence for a unique resource of Nile
751    tilapia in Lake Tanganyika, East Africa. *Environmental Biology of Fishes*, *102*(8), 1107–
752    1117.
753 Syaifudin, M., Bekaert, M., Taggart, J. B., Bartie, K. L., Wehner, S., Palaiokostas, C., Khan, M.
754    G. Q., Selly, S.-L. C., Hulata, G., D'Cotta, H., Baroiller, J.-F., McAndrew, B. J., & Penman,
755    D. J. (2019). Species-Specific Marker Discovery in Tilapia. *Scientific Reports*, Vol. 9.
756 Thodesen, J., Rye, M., Wang, Y.-X., Li, S.-J., Bentsen, H. B., & Gjedrem, T. (2013). Genetic
757    improvement of tilapias in China: Genetic parameters and selection responses in growth,
758    pond survival and cold-water tolerance of blue tilapia (*Oreochromis aureus*) after four
759    generations of multi-trait selection. *Aquaculture* , *396-399*, 32–42.
760 Trewavas, E. (1983). Tilapia Fishes. *British Museum (Natural History), London*, *152*.
761 Waiswa Mwanja, W., Fuerst, P. A., & Kaufman, L. (2012). Reduction of the "ngege",
762    *Oreochromis esculentus* (Teleostei: Cichlidae) populations, and resultant population genetic
763    status in the Lake Victoria Region. *Uganda Journal of Agricultural Sciences*, *13*(2), 65–82.
764 Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer.
765 Wringe, B. F., Anderson, E. C., Jeffery, N. W., Stanley, R. R. E., & Bradbury, I. R. (2019).
766    Development and evaluation of SNP panels for the detection of hybridization between wild
767    and escaped Atlantic salmon (*Salmo salar*) in the western Atlantic. *Canadian Journal of
768    Fisheries and Aquatic Sciences. 76*(5), 695–704.
769 Zhao, H., Fuller, A., Thongda, W., Mohammed, H., Abernathy, J., Beck, B., & Peatman, E.
770    (2019). SNP panel development for genetic management of wild and domesticated white
771    bass (*Morone chrysops*). *Animal Genetics*, *50*(1), 92–96.
772 Zheng, X., Levine, D., Shen, J., Gogarten, S. M., Laurie, C., & Weir, B. S. (2012). A high-
773    performance computing toolset for relatedness and principal component analysis of SNP
774    data. *Bioinformatics* , *28*(24), 3326–3328.