

1 **Title:**

2 Medaka population genome structure and demographic history described via genotyping-by-
3 sequencing

4

5 **Authors and Affiliations:**

6 Takafumi Katsumura^{1,2*}, Shoji Oda³, Mitani Hiroshi⁴, Hiroki Oota^{1*}.

7

8 ¹ Department of Anatomy, Kitasato University School of Medicine

9 1-15-1 Kitasato, Minami-Ku Sagamihara, Kanagawa 252-0374, Japan

10 ² Graduate School of Natural Science and Technology, Okayama University, 1-1-1

11 Tsushimanaka, Kita-Ku, Okayama 700-8530, Japan

12 ³ Department of Integrated Biosciences, Graduate School of Frontier Sciences, University of

13 Tokyo, 5-1-5 Kashiwanoha, Kashiwa, Chiba 277-8562, Japan

14

15 ***Authors for Correspondence:**

16 Takafumi Katsumura, Graduate School of Natural Science and Technology, Okayama

17 University, Okayama, Japan, +81 86 252 7860, tk@okayama-u.ac.jp

18 Hiroki Oota, Department of Anatomy, Kitasato University School of Medicine, Kanagawa,

19 Japan, +81 42 778 9022, hiroki_oota@med.kitasato-u.ac.jp

20 Abstract

21 Medaka is a model organism in medicine, genetics, developmental biology and population
22 genetics. Lab stocks composed of more than 100 local wild populations are available for
23 research in these fields. Thus, medaka represents a potentially excellent bioresource for
24 screening disease-risk- and adaptation-related genes in genome-wide association studies.
25 Although the genetic population structure should be known before performing such an
26 analysis, a comprehensive study on the genome-wide diversity of wild medaka populations
27 has not been performed. Here, we performed genotyping-by-sequencing (GBS) for 81 and
28 12 medakas captured from a bioresource and the wild, respectively. Based on the GBS data,
29 we evaluated the genetic population structure and estimated the demographic parameters
30 using an approximate Bayesian computation (ABC) framework. The autosomal data
31 confirmed that there were substantial differences between local populations and supported our
32 previously proposed hypothesis on medaka dispersal based on mitochondrial genome
33 (mtDNA) data. A new finding was that a local group that was thought to be a hybrid
34 between the northern and the southern Japanese groups was actually a sister group of the
35 northern Japanese group. Thus, this paper presents the first population-genomic study of
36 medaka and reveals its population structure and history based on autosomal diversity.

37

38 Key words: local population, freshwater fish, demography, RAD-seq, bioresource

39 Introduction

40 Medaka (*Oryzias latipes*) is a small fresh-water fish native to East Asia that has attracted
41 attention as a vertebrate model for population genetics (Oota and Mitani 2011; Spivakov *et al.*
42 2014). Wild medaka populations have been maintained in certain universities and research
43 institutes as a bioresource (hereafter, wild lab stocks) with funding from the Japanese
44 government since 1985 (Shima *et al.* 1985). These populations consist of more than 100
45 local-wild populations that have various phenotypic traits (Watanabe-Asaka *et al.* 2014;
46 Igarashi *et al.* 2017) and abundant genetic diversity (Kasahara *et al.* 2007). Geographical
47 features affect the population structure of organisms. Seas and mountains restrict the
48 movement of terrestrial animals and freshwater fish, and the climate also changes the
49 breeding timing and foraging environment. Particularly, noticeable climate differences are
50 observed within island groups, such as the Japanese archipelago, where a large latitudinal
51 difference exists between the southern and northern ends. Therefore, animals are exposed to
52 various selective pressures according to the geographic environment. The local populations
53 have differentiated into local groups by genetic drift, resulting in genetically divergent groups.
54 Medakas stocked as a bioresource are thought to have retained the genetically adapted traits
55 they acquired from various environments of the Japanese archipelago.

56 Our ultimate goal in exploiting medaka characteristics is to establish an experimental
57 system for testing the functional differences between alleles detected by, for instance,
58 genome-wide association studies. Once genetic polymorphisms related to phenotypic traits
59 in medaka are detected, *in vivo* experiments, such as genome editing experiments (Ansai and
60 Kinoshita 2014), can be conducted to understand the functions of the genetic variants
61 (Shimmura *et al.* 2018). Revealing the functional difference between alleles in wild
62 populations allows us to infer the role of genetic polymorphisms in human homologous genes
63 (Igarashi *et al.* 2017; Shimmura *et al.* 2018).

64 Most previous analyses of medaka genetic diversity and population structure have
65 been conducted using mitochondrial DNA (mtDNA). Medaka is divided into four
66 mitochondrial groups: the northern Japanese (N.JPN), southern Japanese (S.JPN), eastern
67 Korean (E.KOR) and western Korean/Chinese groups (W.KOR) (Sakaizumi *et al.* 1980;
68 1983; Sakaizumi 1986; Katsumura *et al.* 2009). These groups do not have sympatric
69 habitats. Although previous allozyme-based studies show that an ambiguous group exists in
70 the geographic boundary region known as Tajima-Tango (see also fig. 1, supplementary fig
71 1), between N.JPN and S.JPN, which is thought to be a hybridization of the two groups
72 (Sakaizumi 1984; Takehana *et al.* 2016), the process has not been fully verified.

73 Based on mtDNA cytochrome *b* gene sequences, N.JPN and S.JPN are composed of 3
74 and 11 subgroups, respectively (Takehana *et al.* 2003). Each subgroup is composed of local
75 populations, and the between-population genetic diversities are greater than the within-
76 population genetic diversities, indicating substantial genetic differentiation between local
77 populations (Katsumura *et al.* 2009). Their habitat environments are also largely different.
78 For instance, there is a large amount of snowfall in the habitats of N.JPN, where the breeding
79 season is short. The habitat of S.JPN is wider than that of N.JPN, and its climate
80 environment is also diverse; e.g., the differences of the annual average temperature and
81 rainfall between the south end (Nago in Okinawa) and north end (Ichinoseki in Iwate)
82 inhabited by S.JPN were 11.6 °C and 988 mm in 2017, respectively
83 (<http://www.jma.go.jp/jma/indexe.html>). The phylogenetic data from these studies also
84 suggested that N.JPN and S.JPN have been spreading in the Japanese archipelago at different
85 times. Particularly, the origin of S.JPN, which has the largest habitat, has been suggested to
86 be the northern part of Kyushu Island based on mtDNA (“Out of Northern Kyushu”
87 hypothesis) (Katsumura *et al.* 2012). However, mtDNA is insufficient to describe the
88 population structure and history because of its single locus. Genetic diversity based on

89 autosomes is essential to understand the medaka population structure, but the comprehensive
90 data of autosomes are still limited.

91 An inference of the population structure and history induced by autosomal information
92 is more robust than that induced by single loci. To unravel the population structure based on
93 autosomes, we comprehensively performed a population-genetic analysis based on autosomal
94 single-nucleotide polymorphisms (SNPs) using all 81 local populations maintained as wild
95 lab stocks at the University of Tokyo. We examined an individual sampled from each
96 population stock and 12 wild individuals captured from the northern part of Kyushu Island,
97 where (Katsumura *et al.* 2012) the medakas currently distributed along the entire Pacific side
98 of the Japanese archipelago originated, to estimate each population's time of expansion from
99 Northern Kyushu. To obtain SNP data, we conducted genotyping-by-sequencing (GBS)
100 (Elshire *et al.* 2011; Narum *et al.* 2013) using a high-throughput sequencer. This method,
101 which allowed us to cost-effectively genotype tens of thousands of SNPs (Andrews *et al.*
102 2016), resulted in an accurate enhancement of the population-genetic estimations because the
103 variance in the demographic parameter estimates decreased when we used many SNPs.
104 Eventually, we obtained more than ten thousand SNPs from eighty-one wild lab stocks and
105 twelve wild-captured medakas. Here, we re-evaluated the medaka genetic diversity and
106 population structures based on these SNPs and reconstructed the population history by
107 assessing three demographic events. These data redefine the medaka local groups and
108 provide a basis of the population history for discussing the role of phenotype-associated
109 alleles in the context of adaptation.

110 Materials and Methods

111 **Samples**

112 We sampled 81 male medakas from 81 wild lab stocks (14 from N.JPN, 56 from S.JPN, 5
113 from Tajima-Tango in the geographic boundary region between N.JPN and S.JPN, 3 from
114 E.KOR, and 3 from W.KOR; note that these groupings come from previous mtDNA
115 sequences and allozyme patterns) at the University of Tokyo in 2014. The lab stocks
116 originated from geographically distinct populations in East Asia and have been maintained
117 since 1985 (Shima *et al.* 1985) as closed colonies in the Graduate School of Frontier Sciences,
118 the University of Tokyo, Kashiwa City (fig. 1). These lab stocks maintain the genetic
119 diversity originating from their habitat (Katsumura *et al.* 2014; Igarashi *et al.* 2017) and show
120 less diversity within populations than between populations (Katsumura *et al.* 2009).
121 Because of these characteristics, we considered that one individual sampled from the lab stock
122 would be adequate to represent its originating population. In addition, we used two wild-
123 captured medaka populations from the Saga Prefecture in the northern part of Kyushu, Japan.
124 One population was the Ogi (SO) population from Southern Saga, and the other population
125 was the Umeshiro (US) population from Northern Saga, both of which were captured in
126 September 2010 (see Katsumura *et al.* 2012). Six of 48 medakas from each wild-captured
127 population were selected randomly and analyzed via GBS.

128

129 **DNA extraction and genotyping-by-sequencing**

130 One-third of the medaka body was dissolved in a 600 μ l lysis buffer containing 1.24% SDS,
131 0.124 M EDTA and 0.062 mg/ml proteinase K (final concentrations). The total genomic
132 DNA was extracted and purified using phenol-chloroform and isopropanol precipitation.
133 After a 70% EtOH wash, an isolated DNA pellet was resuspended in 100 μ l TE buffer and
134 then treated with RNase A (final conc. 1 mg/ml) for 1 hour at room temperature. Then, the

135 DNA was purified again using phenol-chloroform and isopropanol precipitation. For ninety-
136 three samples, the GBS process was outsourced to Macrogen Japan in Kyoto. The
137 procedures for constructing libraries and performing Illumina HiSeq 2000 single-end
138 sequencing were the same as those described by Poland *et al.* (Poland *et al.* 2012) except for
139 the use of the restriction enzyme ApeKI instead of EcoRI–MspI. The sequence lengths were
140 51 bp and included each individual in-line barcode (4–9 bp) for the individual sample. The
141 data have been submitted to the DDBJ Sequence Read Archive (DRA) database under project
142 accession ID: DRA006353.

143

144 **Quality filtering and SNP extraction**

145 Our single-end reads were filtered using *FASTQ Quality Filter* in *FASTX-Toolkit version*
146 *0.0.13* (<http://hannonlab.cshl.edu/fastx-toolkit>) using the following options: -Q 33 -v -z -
147 q 30 -p 90. The draft genome of the medaka sequenced by the PacBio sequencer (Medaka-
148 Hd-rR-pacbio_version2.0.fasta; http://utgenome.org/medaka_v2/#!Assembly.md) was used to
149 align the reads using *BWA backtrack 0.7.12-r1039* (Li and Durbin 2009) using the “-n 0.06”
150 option. After the mapping process, the multi-mapped reads were removed using *Samtools*
151 *v1.2* (Li *et al.* 2009) and the “-Sq 20” option. Following this pipeline, sequencing of the
152 ApeKI-digested GBS libraries generated an average of 3.06 million reads per individual
153 before any quality filtering. The read numbers ranged from 1.89 to 4.10 million reads per
154 individual. After quality filtering, 2.76 million (90.2%) sequences per individual on average
155 were retained, and 0.30 million (9.8%) sequences were eliminated. The retained sequences
156 presented a mean quality score of 38.7 and a GC content of 47.8%. An average of 1.84
157 million of the retained sequences (66.7%) aligned to the medaka autosomal genome, and 0.92
158 million sequences (33.3%) were not aligned and discarded because they mapped to non-
159 autosomal loci (mitochondrial genome and unanchored contigs) and multiple loci.

160 The *Stacks* pipeline (version 1.35) and the *Stacks* workflow
161 (https://github.com/enormandeau/stacks_workflow) were used to generate SNPs and
162 sequences for each individual separately (Catchen *et al.* 2011; 2013). The selected *Stacks*
163 parameters were as follows: minimum stack depth (-m), 3; and number of mismatches when
164 building the catalog (-n), 1. In the ‘rx’ step of *Stacks*, we used the bounded SNP model set
165 to 0.1, the ϵ upper bound set to 0.1 and the log likelihood set to 10. Following this pipeline,
166 five datasets were constructed using the *population* program in *Stacks*.

167 The first was the “PopStat” dataset, which was generated using the “-p 5 -r 0.66”
168 options and mitochondrial grouping (N.JPN, Tajima-Tango, S.JPN, E.KOR, W.KOR; see also
169 fig. 1) to calculate the population-genetic statistics (table 1) using the loci shared across all
170 groups and sequenced in one or more populations in each group.

171 The second was the “Global” dataset, which was generated using the “-p 58 -r 1.00”
172 option and without the mitochondrial grouping, to examine the phylogenetic relationships
173 between the geographic populations and the population structure within the species.

174 The third and the fourth were the “HZ-1” and “HZ-2” datasets (“HZ” is the
175 abbreviation of “Hybrid zone”), respectively, including the 15 boundary populations (N.JPN:
176 Kaga, Maiduru, Miyadu, Obama, and Sabae; Tajima-Tango: Amino, Hamasaka, Kinosaki,
177 Kumihama, and Toyooka; and Honshu: Ayabe, Iwami, Kasumi, Matsue, and Tottori) with or
178 without the Kyushu populations (Kyushu: Fukue, Hiwaki, Izumi, Kadusa, and Kikai). These
179 datasets were generated using the “-p 3 -r 0.70” (without Kyushu) and the “-p 4 -r 0.70” (with
180 Kyushu) option, respectively, to assess the genetic population structure and the history of the
181 boundary population for the Tajima-Tango group.

182 The fifth was the “Local” dataset, including two Kyushu deme samples: one was
183 Umeshiro (US), which was sampled in the northern part of the Saga prefecture, and the other
184 was Ogi (SO), which was sampled in the southern part of the Saga prefecture. These samples

185 were used to estimate the time of Honshu's population divergence from Kyushu to infer the
186 timing of the "Out of Northern Kyushu" event (Katsumura *et al.* 2012).

187 Note that we define the terms "deme samples" and "non-deme samples" as
188 "samples from the local-wild population" and "samples from the wild lab stocks",
189 respectively (see details in Katsumura *et al.* 2009).

190 Against those datasets (except "PopStat" and "HZ-1"), the SNPs with strong linkage
191 disequilibrium were randomly removed, one from each pair of SNPs with $r^2 > 0.2$, using
192 Plink1.9 (Chang *et al.* 2015) and the "--indep-pairwise 12.5 5 0.2 --autosome-num 24" option.
193 These values were set based on a medaka population genomics study (Spivakov *et al.* 2014).
194 Finally, the "Global", "HZ-2" and "Local" datasets contained 8,361 SNPs out of 13,177 SNPs,
195 1,014 SNPs out of 2,661 SNPs and 698 SNPs out of 2,246 SNPs, respectively.

196

197 **Genetic clustering analysis**

198 To obtain a genetic overview of the relationship of medaka geographic populations, we
199 performed a principal component analysis (PCA) of the "Global" dataset as implemented in
200 the *SNPRelate* program (Zheng *et al.* 2012) in R version 3.2.2. Additionally, to examine the
201 genetic relationships between medakas in the Japanese archipelago, we used the subdataset
202 without the E/W.KOR and Chinese populations, which included 7,126 SNPs. We also
203 performed a model-based genetic clustering analysis using *ADMIXTURE v1.23* (Alexander *et al.*
204 *al.* 2009) to estimate the proportions of ancestral medaka populations. We ran 50 replicates
205 with random seeds for the number of clusters (K) from 1 to 9 and calculated the mean of the
206 lowest fivefold cross-validation errors for each K (Jeong *et al.* 2016). Values of K = 4, 5,
207 and 6 showed the first-, second- and third-lowest fivefold cross-validation errors, respectively
208 (supplementary fig. 2). The results of the genetic clustering analyses were visualized by a
209 *ggplot2* package in R (Wickham 2011).

210

211 **Reconstructing the phylogenetic tree using the maximum likelihood method**

212 We considered an individual as representative of a population and generated the individual
213 sequences using the *population* program with the “-p 58 -r 1.00 --phylip_all” option in the
214 *Stacks* pipeline. The dataset included 4,638 partitions and 217,986 bp of nucleotide
215 sequences to compensate for the loci that could not be sequenced by invariable sites across all
216 samples and “N” at the variable site. The dataset including 4,638 partitions and 217,986 bp
217 nucleotide sequences was analyzed via the *IQ-TREE* program (Nguyen *et al.* 2015) to
218 reconstruct a maximum likelihood tree with model selection for each partition (Chernomor *et*
219 *al.* 2016) and 1,000 SH-aLRT/ultrafast-bootstrappings (Guindon *et al.* 2010; Minh *et al.*
220 2013). Then, we used the *FigTree* program (<http://tree.bio.ed.ac.uk/software/figtree/>) to
221 visualize the phylogenetic tree.

222

223 **Inference of demographic parameters of the populations in and around Tajima-Tango** 224 **and the time of “Out of Northern Kyushu” based on the ABC framework**

225 Using the HZ-2 dataset, which added five Kyushu populations to 15 boundary populations of
226 the HZ-1 data set, the population history, including the possibility that the Tajima-Tango
227 group occurred by the admixture of N.JPN and S.JPN, was inferred from 1,014 SNPs using
228 *DIYABC ver2.1.0* (Cornuet *et al.* 2014). We tested four evolutionary scenarios in which
229 N.JPN and S.JPN diverged from an ancestral population at time t_3 : (I) Tajima-Tango
230 originated in N.JPN: Honshu and Kyushu diverged at time t_2 , and then Tajima-Tango
231 diverged from N.JPN at time t_1 (Takehana *et al.* 2016); (II) Admixture of N.JPN and Honshu:
232 Honshu and Kyushu diverged at time t_2 , and then Tajima-Tango occurred at time t_1 by an
233 admixture with rate r between N.JPN and Honshu (Sakaizumi 1984); (III) N.JPN originated in
234 Tajima-Tango: Honshu and Kyushu diverged at time t_2 and then N.JPN diverged from

235 Tajima-Tango at time t_1 ; and (IV) Honshu diverged from Kyushu and then Tajima-Tango
236 occurred by an admixture with Honshu of rate r . For simplicity, all populations were
237 assumed to have constant effective sizes in each lineage (i.e., no bottleneck and expansion).
238 The same prior parameters were defined for four scenarios based on previous studies, and the
239 prior distribution of each parameter is presented in table 2. In addition, we set the
240 conditional constraint as follows: $t_3 > t_2$, $t_3 > t_1$ and $t_2 \geq t_1$. In total, four million
241 simulations were run, which provided approximately one million simulations for each
242 scenario.

243 The summary statistics for all SNP loci included the proportion of loci with null
244 gene diversity, mean gene diversity across polymorphic loci, proportion of loci with null Nei's
245 distance between the two samples, variance across loci of non-null Nei's distances between
246 two samples, proportion of loci with null admixture estimates, mean across loci of non-null
247 admixture estimates, variance across loci of non-null admixture estimates and mean across all
248 locus admixture estimates. The 1% simulated datasets closest to the observed dataset were
249 used to estimate the posterior parameter distributions through a weighted local linear
250 regression procedure (Beaumont *et al.* 2002). Scenarios were compared by estimating their
251 posterior probabilities using the direct estimation and logistic regression methods
252 implemented in *DIYABC* (Cornuet *et al.* 2014). We also estimated the time for "Out of
253 Northern Kyushu" using the Local dataset that included 698 SNPs from *DIYABC*. We
254 inferred the following time-based simple evolutionary scenario: N.Saga diverged from S.Saga
255 at time t because this divergence event was consistent with the "Out of Northern Kyushu"
256 event (Katsumura *et al.* 2012). In this case, all populations were assumed to have
257 unchangeable effective sizes in each lineage.

258 The prior distribution of each parameter is presented in table 2, and the 1%
259 simulated datasets closest to the observed dataset were used to estimate the posterior

260 parameter distributions. We evaluated the accuracy of the demographic parameter
261 estimation by calculating accuracy indicators (*Bias and mean square error* option) in *DIYABC*
262 and accepted the parameters that were non-scaled and scaled by the mean effective population
263 size for the HZ-2 and Local datasets, respectively. Using the infinite-sites model in the
264 Wright-Fisher populations of a constant size, a rough value of N_e was estimated using the
265 relationship $\pi = 4N_e\mu$ (Tajima 1983), where π is genome-wide nucleotide diversity and μ is
266 the mutation rate per site per generation (Osada 2015). We calculated N_e using the above
267 formula with the following values: $\mu = 10^{-9}$ /site/generation; generation = 1 year, which
268 resulted in a mean effective population size from all presented demes in northern Kyushu of
269 1,275,000.

270

271 **Reconstruction of the phylogenetic tree for the partial mitochondrial DNA sequence**

272 To examine the mitochondrial introgression in the tested samples, we generated and analyzed
273 the mitochondrial DNA (mtDNA) partial sequences from the GBS-read data. We mapped
274 the reads to the complete mtDNA sequence of "Clade C" (Matsuda *et al.* 1997; Takehana *et*
275 *al.* 2003), which diverged genetically from S.JPN and N.JPN, to eliminate the mapping bias
276 that occurs when genetically distant sequences are used as a reference. The genomic DNA
277 of "Clade C" was extracted in Katsumura *et al.* 2009. The complete mtDNA sequence was
278 determined as follows. Using PCR, five DNA fragments of mitochondrial DNA were
279 amplified: fragment 1—4,556 bp; fragment 2—4,527 bp; fragment 3-1—4,589 bp; fragment
280 3-2—4,504 bp; and fragment 4—4,546 bp (supplementary fig. 3). Supplementary table 1
281 describes the primers, which were designed based on the inbred strain Hd-rR complete
282 mtDNA sequence (accession number: AP008938). Approximately 20 ng of the genomic
283 DNA was used as a template for the PCR assay in a 50 μ l solution containing dNTP at 0.2
284 mM, 0.2 μ M of each of primer, 0.75 U of EX Taq polymerase HS (TaKaRa Shuzo Co.), and

285 the reaction buffer attached to the polymerase. The reactions were conducted in a TaKaRa
286 PCR Thermal Cycler Dice (TaKaRa Shuzo Co.) using the following protocol: an initial
287 denaturing step at 95°C for 2 min, 40 cycles of denaturation at 95°C for 30 sec, annealing at
288 60°C for 30 sec, extension at 72°C for 300 sec, and a final extension step at 72°C for 5 min.
289 The PCR products were diluted 20-fold and used as templates in the sequencing reaction
290 (following the commercial protocol) with thirty-three primers (supplementary table 1,
291 supplementary fig. 3) and then analyzed in an ABI 3500xL Genetic Analyzer (Life
292 Technologies). The complete mtDNA sequence of "Clade C" was reconstructed using
293 *SeqMan Pro 10.1.2.20* (DNASTAR) and deposited into the international DNA database
294 DDBJ/EMBL/GenBank (accession number: LC335803).

295 The "Clade C" complete mtDNA sequence was used to align the reads using *BWA*
296 *backtrack 0.7.12-r1039* (Li and Durbin 2009) using the "-n 5" option. After the mapping
297 process, to remove the multi-mapped reads, we used *Samtools v1.2* (Li *et al.* 2009) using the
298 "-Sq 20" option. Then, we used *Stacks* with the "-m 3" option for minimum stack depth
299 (Catchen *et al.* 2011; 2013) and then obtained the 322 bp nucleotide sequences. The loci
300 where the sequence was missing were filled in with "N". In addition, the nucleotide position
301 showing the multi-allelic state was replaced by "N". Phylogenetic trees were constructed
302 using the neighbor-joining (NJ) method (Saitou and Nei 1987) with the program MEGA5
303 (Tamura *et al.* 2011). The evolutionary distances were calculated using the Jukes-Cantor
304 method (Jukes and Cantor 1969). The analysis also involved 13 nucleotide sequences from
305 the DNA database (See supplementary fig. 4). All ambiguous positions were removed for
306 each sequence pair. The reliability of the tree was evaluated using 1000 bootstrap replicates
307 (Felsenstein 1985).

308 Results

309 **Autosomal genetic diversity of five groups of *Oryzias latipes***

310 To assess the genetic diversity based on autosomal SNPs within known mitochondrial and
311 hybridization groups—the northern Japanese (N.JPN), southern Japanese (S.JPN), eastern
312 Korean (E.KOR) western Korean/Chinese (W.KOR) and Tajima-Tango groups—their
313 population-genetic summary statistics were calculated using the "PopStat" dataset that
314 included the alignment regions across all groups and contained approximately 45 kb of
315 sequences and 2,453 SNPs (table 1). All groups showed major allele frequencies and
316 heterozygosities of 0.940–0.983 and 0.023–0.054, respectively, indicating many rare alleles
317 and homozygous sites in our wild lab-stocks. W.KOR showed the highest nucleotide
318 diversity (0.0053 ± 0.0003) in all lab-stocks originated from East Asia. S.JPN showed the
319 highest nucleotide diversity (0.0036 ± 0.0001) and N.JPN the lowest in the Japanese
320 archipelago (table 1), which is consistent with the diversity based on mtDNA (Katsumura *et*
321 *al.* 2009).

322

323 **Genetic clustering based on autosomal SNPs inferring medaka population structure**

324 To reveal the medaka population structure based on the autosomal SNPs, we performed
325 genetic clustering analyses using the "Global" dataset. First, we investigated the genetic
326 relationship by a principle component analysis (PCA) using 8,361 SNPs on 24 autosomes.
327 The PCA showed that the wild lab-stocks were divided into two clusters, N.JPN/Tajima-
328 Tango and S.JPN and others from E/W.KOR, which is similar to the grouping based on the
329 mtDNA data (Takehana *et al.* 2003; Katsumura *et al.* 2009)(fig. 2a, supplementary fig. 5).
330 E/W.KOR individuals were plotted dispersedly compared with other clusters, and Yongcheon
331 (E.KOR) was close to N.JPN/Tajima-Tango. As PC3 was divided into E.KOR and W.KOR,
332 the autosomes in these mtDNA groups were differentiated into each other.

333 Although Tajima-Tango was considered to be a hybridization group between N.JPN
334 and S.JPN (Sakaizumi 1984), all individuals from Tajima-Tango were plotted with those from
335 N.JPN as a cluster (fig. 2a). Next, for a fine-scale mapping of Japanese medaka, we
336 performed the PCA based on 7,126 SNPs, excluding E/W.KOR individuals (fig. 2b). The
337 PCA showed that the Kyushu populations (Kudamatsu, Ogi, Izumi, Hiwaki, Kikai, Nago,
338 Gushikami, Hisayama, Umeshiro, Ashibe, Arita, Kusu and Nobeoka) in S.JPN were dispersed
339 along the PC2 axis (supplementary fig. 6), while that Tajima-Tango overlapped with N.JPN
340 again, even though we used up to PC5 (supplementary fig. 5). This result indicated that
341 Tajima-Tango was not considerably differentiated from N.JPN on autosomes, suggesting that
342 S.JPN was differentiated within the groups and the genetic diversity was especially high
343 among the Kyushu populations (fig. 2b).

344 To examine the phylogenetic relationship among 83 populations, we constructed a
345 maximum likelihood (ML) tree based on 4,638 short fragment sequences (comparative
346 nucleotide sequence length: 217,986 bp) using *IQ-TREE*. The maximum likelihood tree
347 showed relatively high bootstrap values on each branch and almost the same topology as
348 previous trees based on mtDNA (Takehana *et al.* 2003; Katsumura *et al.* 2009), except for
349 Yongcheon (fig. 3 and supplementary fig. 7; see also Discussion). Medakas from the Korean
350 peninsula were genetically close to each other in both groups, whereas those from Shanghai,
351 which have been classified into W.KOR based on mtDNA studies, diverged from the other
352 individuals from W.KOR. The medakas in the Japanese archipelago were divided into two
353 major clusters similar to mtDNA. One was the N.JPN/Tajima-Tango, and the other was
354 S.JPN. The N.JPN/Tajima-Tango cluster was further divided into two submajor clusters
355 with a 100% bootstrapping value in the ML tree based on concatenated sequences obtained
356 from GBS, although N.JPN and Tajima-Tango overlapped in PCA, which reduced SNP
357 information, i.e., N.JPN and Tajima-Tango diverged on the whole autosomal sequenced

358 regions. S.JPN was also further divided into two sub-major clusters: Kyushu-only and
359 Kyushu & others, which included 14 Kyushu populations, represented in fig. 3 by an open
360 and a closed red circle, respectively. While the Kyushu-only cluster diverged into the
361 northern and the southern Kyushu, the Kyushu & others cluster diverged to the Pacific side of
362 eastern and the northern Honshu (main-island Japan), which included several clusters
363 supported by high bootstrapping value (SH-aLRT \geq 80% and UFboot \geq 95%, supplementary
364 fig7). Finally, we characterized the S.JPN ancestry in the context of East Asian genetic
365 diversity of medaka by performing an *ADMIXTURE* analysis, which is a model-based
366 unsupervised genetic clustering method. With the optimal number of ancestral components
367 ($K = 4$), S.JPN medakas were assigned to two distinct ancestries (fig. 3 and supplementary
368 fig. 2). In suboptimal runs with more ancestral components ($K = 5, 6$), only S.JPN in
369 Honshu (main-island Japan) was assigned to the other ancestral components found in the
370 northern Kyushu populations. This analysis indicated that genetic diversities in the northern
371 Kyushu populations were the highest among S.JPN. Thus, the genetic clustering analyses
372 based on genome-wide SNP data strongly suggested that S.JPN spread from northern Kyushu
373 to Honshu.

374

375 **Boundary population genomes similar to that of the northern Japanese group**

376 To explore hybridization signatures on autosomes, we examined allele-sharing between
377 boundary populations and surrounding populations using the “HZ-1” dataset (see Materials
378 and Methods). We classified the fixed alleles between the groups into three states: shared by
379 N.JPN and Tajima-Tango, shared by S.JPN and Tajima-Tango, and shared by N.JPN and
380 S.JPN. Then, we summarized each state’s total numbers, as shown in table in fig. 4. We
381 found 1,380 out of 4,661 SNPs that were common alleles between two of the three groups
382 (fig. 4). Regarding those alleles, the majority (81.4%) were common alleles between

383 Tajima-Tango and N.JPN, which was a much higher frequency than that between Tajima-
384 Tango and S.JPN (11.2%). The rest of them were common alleles between N.JPN and
385 S.JPN (7.5%), i.e., specific alleles in Tajima-Tango. These proportions are near those of a
386 previous study based on 96 genomic regions (4 loci for each chromosome) (Takehana *et al.*
387 2016). Although the regions with neighboring common alleles between Tajima-Tango and
388 S.JPN were observed for certain chromosomes, alleles on the Tajima-Tango genome were
389 mostly shared with N.JPN, suggesting that Tajima-Tango is a subgroup of N.JPN with a short
390 divergence time.

391 To investigate the mitochondrial introgression in boundary populations, we
392 reconstructed the phylogenetic tree based on partial mtDNA sequences generated by mapping
393 short reads from the GBS of the complete mitochondrial genome (supplementary fig. 4).
394 The mtDNA sequences from Toyooka and Kinosaki in Tajima-Tango were clustered and
395 closely related to those from Kaga, which was a root population in N.JPN. Considering that
396 Tajima-Tango was a subgroup of N.JPN, which we inferred from the ML tree based on
397 autosomal sequences, these phylogenetic positions on the mtDNA tree would have also
398 reflected an evolutionary history in which N.JPN diverged from the Tajima-Tango group.
399 We confirmed that Amino, Kumihama and Hamasaka in Tajima-Tango had the mitochondrial
400 genomes of phylogenetically distant populations classified into S.JPN (supplementary fig. 4).
401 Additionally, the mitochondrial genome of Ayabe in S.JPN was classified into N.JPN but
402 diverged from Tajima-Tango. These results suggest that mitochondrial genome
403 introgressions occurred reciprocally, meaning that they occurred not only from S.JPN to
404 Tajima-Tango but also from N.JPN to S.JPN.

405

406 **Inferring medaka demographic parameters**

407 To infer the demographic parameters, effective population size (N_e), time (T), and proportion

408 of the admix (r), we analyzed the "HZ-2" dataset based on the coalescent theory using an
409 approximate Bayesian computation (ABC) framework. We performed the model selection
410 to identify the best explanation scenario for the observed data from the four scenarios; (I)
411 Tajima-Tango originated in N.JPN, (II) Admixture of N.JPN and Honshu, (III) N.JPN
412 originated in Tajima-Tango, and (IV) Admixture of Tajima-Tango and Honshu (fig. 5a).
413 *DIYABC* has two different approaches (directional and logistic) for model selection. Based
414 on each criterion, scenario III (Honshu diverged from Kyushu and then N.JPN diverged from
415 the Tajima-Tango group without admixture with Honshu) and scenario IV (Honshu diverged
416 from Kyushu and then Tajima-Tango occurred by an admixture with Honshu at r rate) were
417 supported by the directional (posterior probability: 0.4720, 95% CI: 0.0344–0.9096) and
418 logistic (posterior probability: 0.4477, 95% CI: 0.4313–0.4641) approaches, respectively
419 (supplementary fig. 8). Our data could not distinguish between scenarios III and IV (see the
420 Discussion section).

421 The posterior parameter estimates of scenarios III and IV shown in table 2 and figs.
422 5c and 5d were not scaled because various measures of accuracy (RRMISE, RMeanAD, and
423 RRMSE in the *DIYABC* output) indicated that non-scaled parameters fit the observed data
424 better than scaled parameters. Additionally, although we inferred the time to the most
425 common ancestor at three divergence events (fig. 5), these times might be overestimated
426 because the dataset for this estimation consisted of non-deme samples, i.e., polymorphisms
427 within populations were underestimated, and the random-mating assumption could not be
428 satisfied. Therefore, we estimated the timing of the "Out of Northern Kyushu" event using
429 the "Local" dataset composed of the two deme samples (S.Saga and N.Saga), which were split
430 at first on the S.JPN lineage (fig. 3). Based on a simple hypothesis, constant population size
431 and no migration, we obtained the estimated time (510,000 years ago, 95% CI: 337,875–
432 679,575) for the ancestral divergence of the two deme samples (table 2 and fig. 5e), which

433 was calculated from a scaled parameter by the mean effective population size because various
434 accuracy measures indicated that scaled parameters fit the observed data better than non-
435 scaled parameters. Thus, our population-genetic estimate using ABC suggests that medakas
436 diverged and dispersed throughout Pacific-side Japan approximately 510 kyr ago.

437 Discussion

438 **Redefined subgroups in Japanese archipelago based on the population structure**
439 **inferred from genome-wide single-nucleotide polymorphisms**

440 *Oryzias latipes* can be divided into five groups by mtDNA sequences and allozymic
441 electrophoresis patterns (Sakaizumi 1984; Takehana *et al.* 2003). In this study, based on
442 autosomal SNPs, the genetic clustering analysis showed that "K = 4" was the most supportive
443 because it presented the lowest fivefold cross-validation error, indicating that N.JPN and
444 S.JPN were divided into three ancestral clusters. When the K values increased, only S.JPN
445 divided into more subgroups, which suggests that the S.JPN group was composed of more
446 divergent groups than the other groups. Considering together with the results of the ML tree
447 analysis, it is possible to redefine subgroups composed of each major group for our wild lab-
448 stocks originated from the Japanese archipelago as follows. First, Tajima-Tango, which had
449 been considered a hybridization group, should be included under the N.JPN group because it
450 shows almost the same ancestral component as N.JPN. Then, the group called N.JPN should
451 be assigned to a subgroup, for which we propose the name "Sea of Japan side of Northeastern
452 Honshu (SJ.NEH)." These two subgroups, Tajima-Tango and SJ.NEH, compose the N.JPN
453 group (supplementary fig. 7). Second, S.JPN can be divided into several subgroups, San-in,
454 San-yo/Shikoku/Kinki and the Pacific Ocean side of Northeastern Honshu (PO.NEH),
455 because the Kyushu & others cluster was subdivided into three sub-clusters composed of
456 geographically neighboring populations. Adding the Kyushu subgroup, S.JPN is composed
457 of four subgroups. Thus, medaka in Japanese archipelago could also be composed of the six
458 distinct subgroups based on autosomal genetic diversity (supplementary fig. 7).

459 S.JPN can be divided into finer subgroups based on the mitochondrial genome
460 (Takehana *et al.* 2003), likely because the effective population size of mtDNA is one quarter
461 that of a nuclear gene. This causes the intermingled branch pattern of the mtDNA tree,

462 which is not associated with geographic distance. Although this pattern has been observed
463 our wild lab-stocks in a previous study based on mtDNA (Katsumura *et al.* 2009), it has been
464 interpreted as artificial migrations accompanied by recent human activities (Takehana *et al.*
465 2003). If the branching pattern was formed by human activities, various ancestral
466 components should appear independent of geographic relatedness in the result of the
467 *ADMIXTURE* analysis based on autosomal SNPs. However, no signal of recent migration
468 was found in the *ADMIXTURE* result. Rather, the autosomal SNP data support another
469 hypothetical scenario in which the mtDNA tree topology reflects ancestral polymorphisms
470 only and their local fixation is caused by a small effective population size (Katsumura *et al.*
471 2009). This result indicates that arguing the geographical origins of medaka based only on
472 mtDNA may lead to false conclusions.

473 From our data, it may be difficult to accurately evaluate the genetic diversities in
474 E/W.KOR groups because the number of populations examined in each group was small.
475 Therefore, though care must be taken in the interpretation, W.KOR showed the highest
476 nucleotide diversity and E.KOR the second highest among the four major mitochondrial
477 groups (table 1). The W.KOR group included the Chinese medaka in Shanghai, which could
478 have elevated its value, while the latter group did not include any geographically distant
479 populations. The clustering analysis showed that Shanghai from W.KOR and Yongcheon
480 from E.KOR had an ancestry component from N.JPN (fig.3). This suggests two
481 possibilities: one is that the N.JPN ancestor was derived from the E/W.KOR ancestor, and the
482 other is that contamination occurred through the maintained wild lab stocks. To investigate
483 these possibilities and the origin of *O. latipes*, a population-based genome wide analysis must
484 be conducted to increase the population numbers of E/W.KOR and include the sister species
485 *O. curvinotus* and *O. luzonensis*.

486

487 **Reconstruction of the medaka population history in the Japanese archipelago**

488 Our genome-wide analysis shows that medakas in N.JPN and S.JPN are deeply divergent and
489 dispersed over the Japanese archipelago from different locations at different times. In
490 particular, our ABC analysis indicates that SJ.NEH originated in and diverged from Tajima-
491 Tango. This inference about the history of N.JPN after divergence from S.JPN is not
492 consistent with previous inferences from the allozyme, mtDNA and limited autosomal SNP
493 analyses, which have suggested that Tajima-Tango is a hybridization group between PO.NEH
494 and S.JPN, or Tajima-Tango is a sub-group derived from PO.NEH. Our GBS data show that
495 (i) the nucleotide diversity in Tajima-Tango is higher than that in SJ.NEH (table 1), (ii) S.JPN
496 is genetically more closely related to Tajima-Tango than to SJ.NEH based on the allele-
497 sharing rates (fig. 4), and (iii) the Tajima-Tango branch is the root in the N.JPN clade of the
498 phylogenetic tree based on partial mtDNA sequences (supplementary fig. 4), which is also
499 shown by the whole mitochondrial genome analysis (Hirayama *et al.* 2010). These data
500 suggest that Tajima-Tango forms an outgroup to all present-day SJ.NEH and it spread along
501 the Sea of Japan side. Furthermore, the ABC framework's estimation supports our scenario,
502 although the analysis does not statistically distinguish between scenario III and scenario IV.
503 Even if the admixture occurred, the inferred ratio of the admixture is too low (fig. 5d).

504 Our findings strongly support the “Out of northern Kyushu” model of S.JPN
505 proposed in Katsumura *et al.* (2012) and have revealed the dispersal route of SJ.NEH/Tajima-
506 Tango in N.JPN. The genetic clustering analyses and phylogenetic tree based on the GBS
507 data elucidate medaka history better than previous mitochondrial DNA analyses. In
508 particular, the *ADMIXTURE* analysis shows that the two ancestry components (yellow and
509 blue in fig. 3) were observed in northern Kyushu, suggesting that S.JPN dispersed in three
510 different directions after dispersing out of Northern Kyushu. Geographically, Shikoku, San-
511 yo, and Kinki are separated by sea, but medakas in the three local lands share the same

512 ancestral component (blue in fig.6; see also supplementary fig. 1 for geographic information).
513 Medaka is a freshwater fish, but survival and reproduction are possible even in seawater
514 (Inoue and Takei 2002). Because most of the rivers in the Japanese archipelago are steep
515 and short, the rivers tend to flood after heavy rain. Although medakas are highly likely to
516 drift to the sea on each occasion, medakas can have survived even in the sea and may have
517 returned to the river because of their saltwater tolerance. Thus, the possibility of moving
518 from river to river through the sea cannot be ignored. The results of this study suggest a
519 history in which medakas migrated throughout the Japanese archipelago through the sea.

520 The different branching patterns of autosomal and mtDNA trees suggest that the
521 mtDNA introgression occurred not only from S.JPN to Tajima-Tango but also from N.JPN to
522 S.JPN. The Tajima-Tango region is surrounded by mountains (supplementary fig. 1);
523 however, it contains the lowest watershed (sea level 95.45 m) in the Japanese archipelago.
524 The medaka has possibly moved in both directions across the watershed. From the above,
525 the most plausible scenario is as follows (fig. 6). Ancestral S.JPN and N.JPN diverged first
526 and independently reached their current habitats in the Japanese archipelago. After S.JPN,
527 which is an ancestor of San-in, San-yo/Shikoku/Kinki and PO.NEH, dispersed from northern
528 Kyushu approximately 510 kyr ago, SJ.NEH diverged from Tajima-Tango in N.JPN. Then,
529 their descendant populations spread rapidly to northern Honshu on the Sea of Japan side.
530 Meanwhile, S.JPN dispersed in as many as three different directions and then spread rapidly
531 northeastward from the western part of Fossa Magna. In the process of dispersing across the
532 main island of Japan, certain S.JPN populations infiltrated the Tajima-Tango region from the
533 west and the south, and the resulting mtDNA introgression occurred independently.

534

535 Conclusion

536 Our genome-wide SNP analysis reconstructed the detailed population structure and reliable

537 history of medaka that evolved in the Japanese archipelago. Since the distribution of the
538 subgroups was highly consistent with the geographical features, several adaptive traits could
539 have evolved in each subgroup. Furthermore, the boundary populations were not caused by
540 a hybridization event but instead were the origin of the populations dispersed to a
541 northeastern part of the Japanese archipelago on the Sea of Japan side. A better
542 understanding of the population structure and history of medaka will support association
543 studies for phenotypes and genotypes related to environmental adaptation.
544

545 References

- 546 Alexander, D. H., J. Novembre, and K. Lange, 2009 Fast model-based estimation of ancestry
547 in unrelated individuals. *Genome Research* 19: 1655–1664.
- 548 Andrews, K. R., J. M. Good, M. R. Miller, G. Luikart, and P. A. Hohenlohe, 2016 Harnessing
549 the power of RADseq for ecological and evolutionary genomics. *Nature Review Genetics*
550 17: 81–92.
- 551 Ansai, S., and M. Kinoshita, 2014 Targeted mutagenesis using CRISPR/Cas system in
552 medaka. *Biology Open* 3: 362–371.
- 553 Beaumont, M. A., W. Zhang, and D. J. Balding, 2002 Approximate Bayesian computation in
554 population genetics. *Genetics* 162: 2025–2035.
- 555 Catchen, J., P. A. Hohenlohe, S. Bassham, A. Amores, and W. A. Cresko, 2013 Stacks: an
556 analysis tool set for population genomics. *Mol Ecol* 22: 3124–3140.
- 557 Catchen, J. M., A. Amores, P. Hohenlohe, W. Cresko, and J. H. Postlethwait, 2011 Stacks:
558 building and genotyping Loci de novo from short-read sequences. *G3 (Bethesda)* 1: 171–
559 182.
- 560 Chang, C. C., C. C. Chow, L. C. Tellier, S. Vattikuti, S. M. Purcell *et al.*, 2015 Second-
561 generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* 4: 7.
- 562 Chernomor, O., A. von Haeseler, and B. Q. Minh, 2016 Terrace Aware Data Structure for
563 Phylogenomic Inference from Supermatrices. *Systematic Biology* 65: 997–1008.
- 564 Cornuet, J.-M., P. Pudlo, J. Veyssier, A. Dehne-Garcia, M. Gautier *et al.*, 2014 DIYABC
565 v2.0: a software to make approximate Bayesian computation inferences about population
566 history using single nucleotide polymorphism, DNA sequence and microsatellite data.
567 *Bioinformatics* 30: 1187–1189.
- 568 Elshire, R. J., J. C. Glaubitz, Q. Sun, J. A. Poland, K. Kawamoto *et al.*, 2011 A robust, simple
569 genotyping-by-sequencing (GBS) approach for high diversity species. *PLoS ONE* 6:
570 e19379.
- 571 Felsenstein, J., 1985 CONFIDENCE LIMITS ON PHYLOGENIES: AN APPROACH
572 USING THE BOOTSTRAP. *Evolution* 39: 783–791.
- 573 Guindon, S., J.-F. Dufayard, V. Lefort, M. Anisimova, W. Hordijk *et al.*, 2010 New
574 algorithms and methods to estimate maximum-likelihood phylogenies: assessing the
575 performance of PhyML 3.0. *Systematic Biology* 59: 307–321.
- 576 Hirayama, M., T. Mukai, M. Miya, Y. Murata, Y. SEKIYA *et al.*, 2010 Intraspecific variation
577 in the mitochondrial genome among local populations of Medaka *Oryzias latipes*. *Gene*
578 457: 13–24.
- 579 Igarashi, K., J. Kobayashi, T. Katsumura, Y. Urushihara, K. Hida *et al.*, 2017 An Approach to
580 Elucidate NBS1 Function in DNA Repair Using Frequent Nonsynonymous
581 Polymorphism in Wild Medaka (*Oryzias latipes*) Populations (S. D. Fugmann, Ed.). *PLoS*
582 *ONE* 12: e0170006–19.

- 583 Inoue, K., and Y. Takei, 2002 Diverse adaptability in oryzias species to high environmental
584 salinity. *Zool. Sci.* 19: 727–734.
- 585 Jeong, C., S. Nakagome, and A. Di Rienzo, 2016 Deep History of East Asian Populations
586 Revealed Through Genetic Analysis of the Ainu. *Genetics* 202: 261–272.
- 587 Jukes, T. H., and C. R. Cantor, 1969 Evolution of protein molecules. *Mammalian protein*
588 *metabolism* 3: 132.
- 589 Kasahara, M., K. Naruse, S. Sasaki, Y. Nakatani, W. Qu *et al.*, 2007 The medaka draft
590 genome and insights into vertebrate genome evolution. *Nature* 447: 714–719.
- 591 Katsumura, T., S. Oda, S. Mano, N. Suguro, K. Watanabe *et al.*, 2009 Genetic differentiation
592 among local populations of medaka fish (*Oryzias latipes*) evaluated through grid- and
593 deme-based sampling. *Gene* 443: 170–177.
- 594 Katsumura, T., S. Oda, S. Nakagome, T. Hanihara, H. Kataoka *et al.*, 2014 Natural allelic
595 variations of xenobiotic-metabolizing enzymes affect sexual dimorphism in *Oryzias*
596 *latipes*. *Proceedings of the Royal Society B: Biological Sciences* 281:.
- 597 Katsumura, T., S. Oda, K. Tsukamoto, T. Yamashita, M. Aso *et al.*, 2012 A population
598 genetic study on the relationship between medaka fish and the spread of wet-rice
599 cultivation across the Japanese archipelago. *Anthropol. Sci.*
- 600 Li, H., and R. Durbin, 2009 Fast and accurate short read alignment with Burrows-Wheeler
601 transform. *Bioinformatics* 25: 1754–1760.
- 602 Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan *et al.*, 2009 The Sequence
603 Alignment/Map format and SAMtools. *Bioinformatics* 25: 2078–2079.
- 604 Matsuda, M., H. Yonekawa, S. Hamaguchi, and M. Sakaizumi, 1997 Geographic variation
605 and diversity in the mitochondrial DNA of the medaka, *Oryzias latipes*, as determined by
606 restriction endonuclease analysis. *Zool. Sci.* 14: 517–526.
- 607 Minh, B. Q., M. A. T. Nguyen, and A. von Haeseler, 2013 Ultrafast approximation for
608 phylogenetic bootstrap. *Molecular Biology and Evolution* 30: 1188–1195.
- 609 Narum, S. R., C. A. Buerkle, J. W. Davey, M. R. Miller, and P. A. Hohenlohe, 2013
610 Genotyping-by-sequencing in ecological and conservation genomics. *Mol Ecol* 22: 2841–
611 2847.
- 612 Nguyen, L.-T., H. A. Schmidt, A. von Haeseler, and B. Q. Minh, 2015 IQ-TREE: a fast and
613 effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular*
614 *Biology and Evolution* 32: 268–274.
- 615 Oota, H., and H. Mitani, 2011 Human Population Genetics Meets Medaka, pp. 339–350 in
616 *Medaka: A Model for Organogenesis, Human Disease, and Evolution*, edited by K.
617 Naruse, M. Tanaka, and H. Takeda. *Medaka: A Model for Organogenesis, Human*
618 *Disease, and Evolution*, Springer Japan, Tokyo.
- 619 Osada, N., 2015 Genetic diversity in humans and non-human primates and its evolutionary
620 consequences. *Genes Genet. Syst.* 90: 133–145.

- 621 Poland, J. A., P. J. Brown, M. E. Sorrells, and J.-L. Jannink, 2012 Development of High-
622 Density Genetic Maps for Barley and Wheat Using a Novel Two-Enzyme Genotyping-
623 by-Sequencing Approach (T. Yin, Ed.). PLoS ONE 7: e32253–8.
- 624 Saitou, N., and M. Nei, 1987 The neighbor-joining method: a new method for reconstructing
625 phylogenetic trees. *Molecular Biology and Evolution* 4: 406–425.
- 626 Sakaizumi, M., 1986 Genetic divergence in wild populations of Medaka, *Oryzias latipes*
627 (Pisces: Oryziatidae) from Japan and China. *Genetica*.
- 628 Sakaizumi, M., 1984 Rigid Isolation between the Northern Population and the Southern
629 Population of the Medaka, *Oryzias latipes* (Genetics). *Zool. Sci.* 1: 795–800.
- 630 Sakaizumi, M., N. EGAMI, and K. MORIWAKI, 1980 Allozymic Variation in Wild
631 Populations of the Fish, *Oryzias latipes*. *Proc. Jpn. Acad., Ser. B* 56: 448–451.
- 632 Sakaizumi, M., K. Moriwaki, and N. Egami, 1983 Allozymic variation and regional
633 differentiation in wild populations of the fish *Oryzias latipes*. *Copeia* 1983: 311.
- 634 Shima, A., A. Shimada, M. Sakaizumi, and N. Egami, 1985 First listing of wild stocks of the
635 medaka *Oryzias latipes* currently kept by the Zoological Institute, Faculty of Science,
636 University of Tokyo. *Journal of the Faculty of Science, Imperial University of Tokyo*
637 Sect IV, *Zoology* 16: 27–35.
- 638 Shimmura, T., T. Nakayama, A. Shinomiya, S. Fukamachi, M. Yasugi *et al.*, 2018 Dynamic
639 plasticity in phototransduction regulates seasonal changes in color perception. *Nature*
640 *Communications* 1–7.
- 641 Spivakov, M., T. O. Auer, R. Peravali, I. Dunham, D. Dolle *et al.*, 2014 Genomic and
642 phenotypic characterization of a wild medaka population: towards the establishment of an
643 isogenic population genetic resource in fish. 4: 433–445.
- 644 Tajima, F., 1983 Evolutionary relationship of DNA sequences in finite populations. *Genetics*
645 105: 437–460.
- 646 Takehana, Y., N. Nagai, M. Matsuda, K. Tsuchiya, and M. Sakaizumi, 2003 Geographic
647 variation and diversity of the cytochrome b gene in Japanese wild populations of medaka,
648 *Oryzias latipes*. *Zool. Sci.* 20: 1279–1291.
- 649 Takehana, Y., M. Sakai, T. Narita, T. Sato, K. Naruse *et al.*, 2016 Origin of Boundary
650 Populations in Medaka (*Oryzias latipes* Species Complex). *Zool. Sci.* 33: 125–131.
- 651 Tamura, K., D. Peterson, N. Peterson, G. Stecher, M. Nei *et al.*, 2011 MEGA5: Molecular
652 Evolutionary Genetics Analysis Using Maximum Likelihood, Evolutionary Distance, and
653 Maximum Parsimony Methods. *Molecular Biology and Evolution* 28: 2731–2739.
- 654 Watanabe-Asaka, T., Y. SEKIYA, H. Wada, T. Yasuda, I. Okubo *et al.*, 2014 Regular
655 heartbeat rhythm at the heartbeat initiation stage is essential for normal cardiogenesis at
656 low temperature. *BMC Dev. Biol.* 14:.
- 657 Wickham, H., 2011 ggplot2. *Wiley Interdisciplinary Reviews: Computational Statistics* 3:
658 180–185.

659 Zheng, X., D. Levine, J. Shen, S. M. Gogarten, C. Laurie *et al.*, 2012 A high-performance
660 computing toolset for relatedness and principal component analysis of SNP data.
661 *Bioinformatics* 28: 3326–3328.

662

663 Acknowledgments

664 We thank Mrs. Shizuko Chiba, Mrs. Sumiko Tomizuka, Dr. Atsuko Shimada and Prof.
665 Emeritus Akihiro Shima (the University of Tokyo) for maintaining the medaka stocks from
666 wild populations. This study was supported in part by grants-in-aid for Young Scientists (B),
667 no. 16K21352 to T.K., and by a grant-in-aid for Scientific Research (A), nos. 26251050,
668 25251046 and 17H01453 to H.O. T.K. was also supported by a grant-in-aid from JSPS
669 Research Fellowships (16J07227).

670 Figure legends

671 **FIG. 1: Map of the original locations of the wild lab stocks and wild-captured medakas.**

672 In the upper central map, the enlarged red frame shows the boundary region between S.JPN
673 and N.JPN. Each color represents the mtDNA and allozyme-based groups shown in the left
674 upper inset box. The numbers on the map are consistent with the population IDs, with the
675 names on the right bottom inset box.

676

677 **FIG. 2: Results of principle component analysis (PCA) using SNPs in East Asia (a) and**

678 **the Japanese archipelago (b).** Each plot shows PC1 versus PC2. The population names
679 of each point are described in fig. S6.

680

681 **FIG. 3: Phylogenetic tree using the maximum likelihood method and an ancestry**

682 **barplot with ADMIXTURE analysis.** Red closed and open circles represent the northern
683 and southern Kyushu populations, respectively. “Taj-Tan” in the tree is the abbreviation for
684 Tajima-Tango.

685

686 **FIG. 4: Shared allele distribution among boundary populations.** The gray rectangle

687 represents the total length of each medaka chromosome. The table in the figure shows the
688 number of alleles shared between the groups observed on each chromosome. Mark
689 represents the state whose group shares that allele. *1 and *2 are the total number from
690 chromosomes 1 to 12 and from chromosomes 13 to 24, respectively.

691

692 **FIG. 5: Scenarios for estimation of the demographic parameters.** Possible scenarios for

693 the population history of the Tajima-Tango group (a) and the "Out of northern Kyushu (NK)"
694 event (b) are shown. The TMRCA estimated using ABC are in scenario3 (c), scenario4 (d)

695 and Out of NK (e).

696

697 **FIG. 6: Map representing the ancestry proportions from ADMIXTURE analysis at $K =$**

698 **6.** Solid and dashed lines represent the spreading patterns of S.JPN and N.JPN inferred by

699 GBS data, respectively.

Figure 1

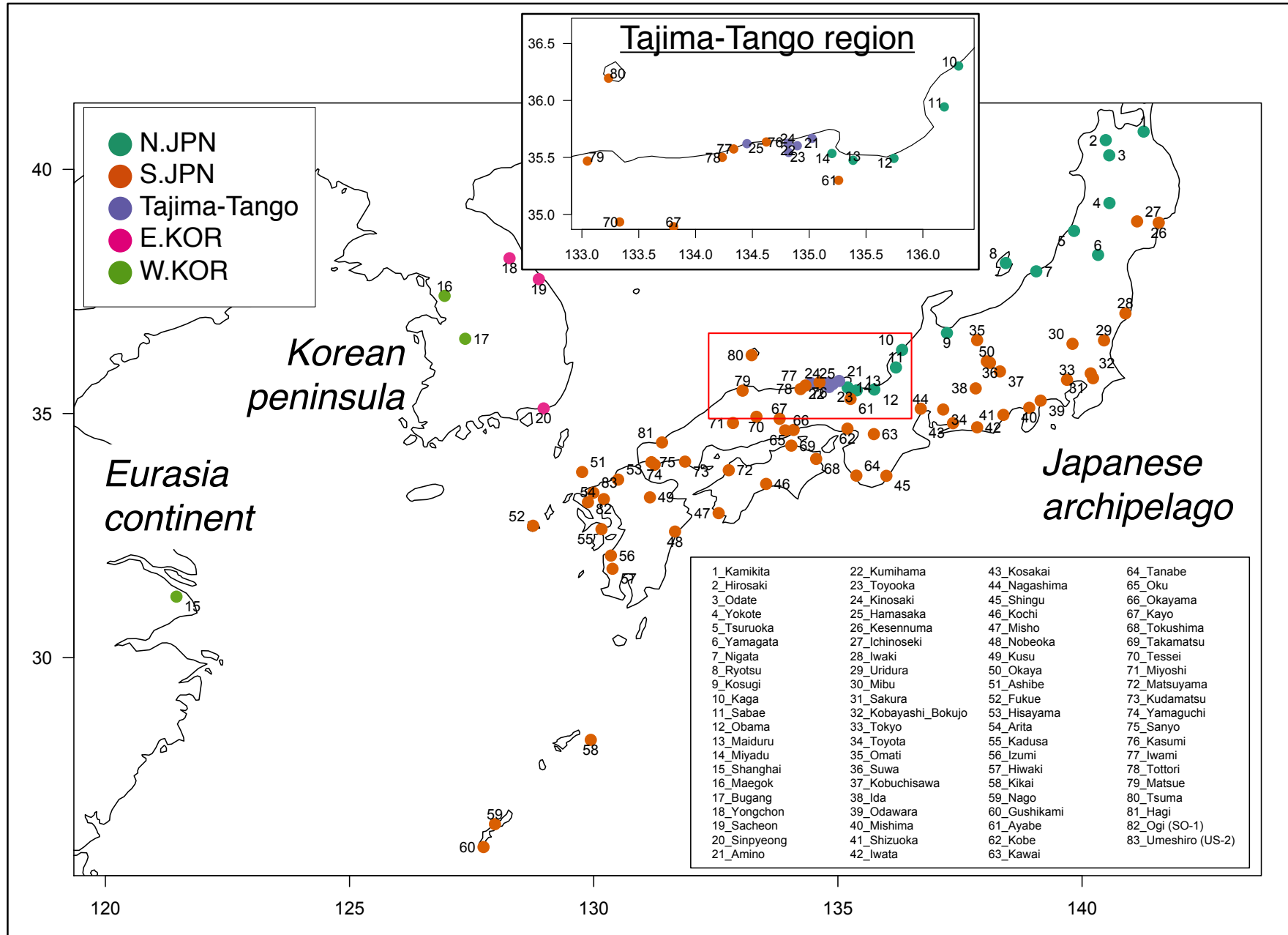


Figure 2

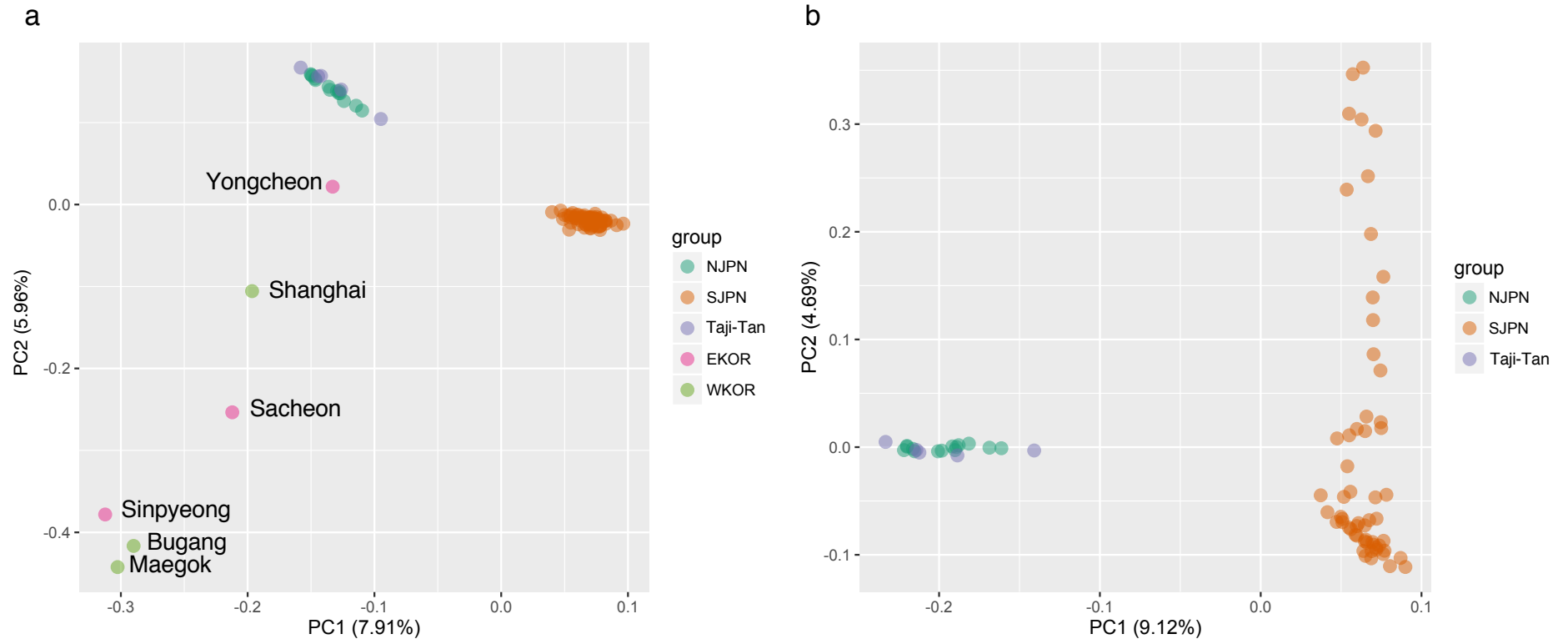


Figure 3

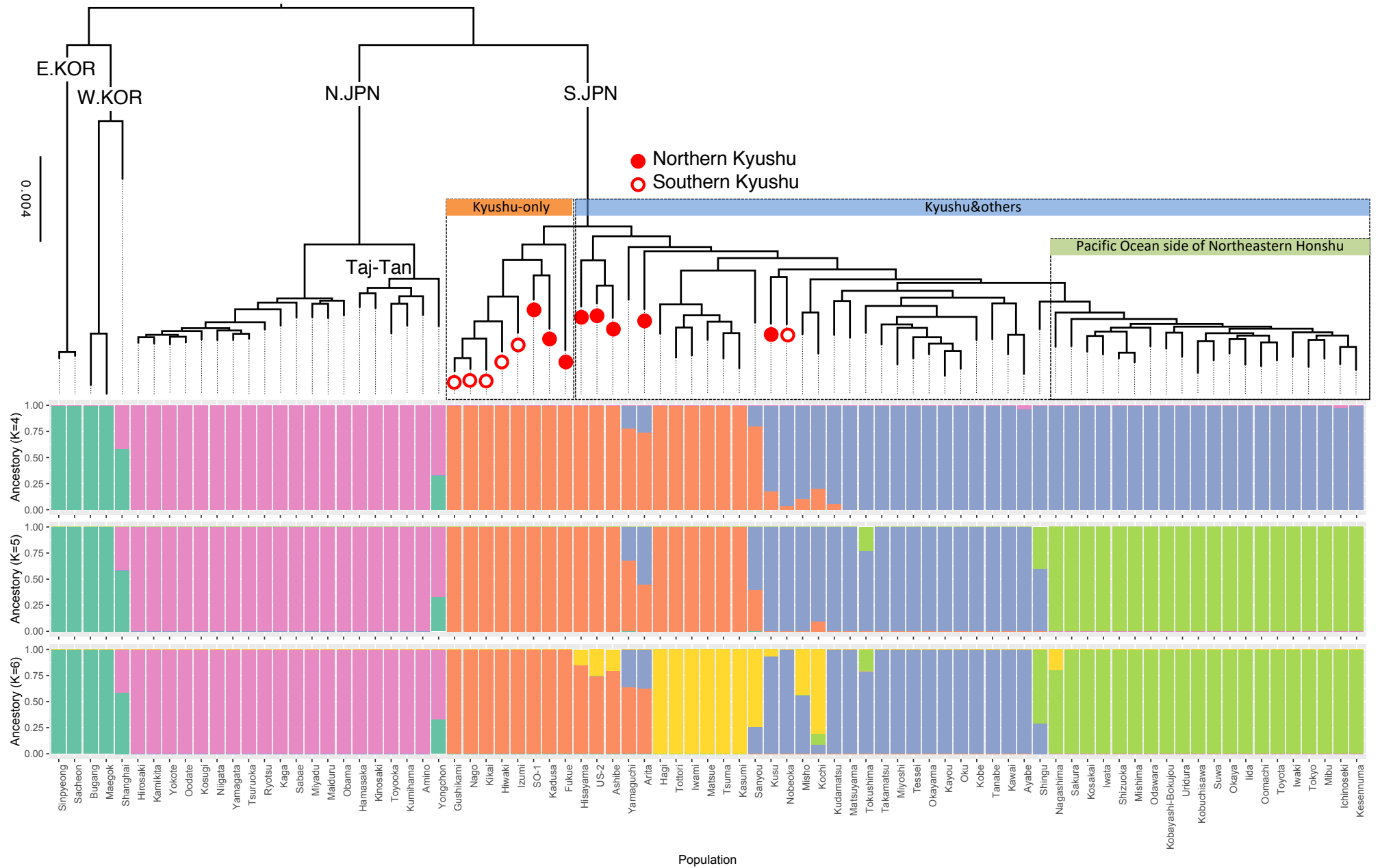
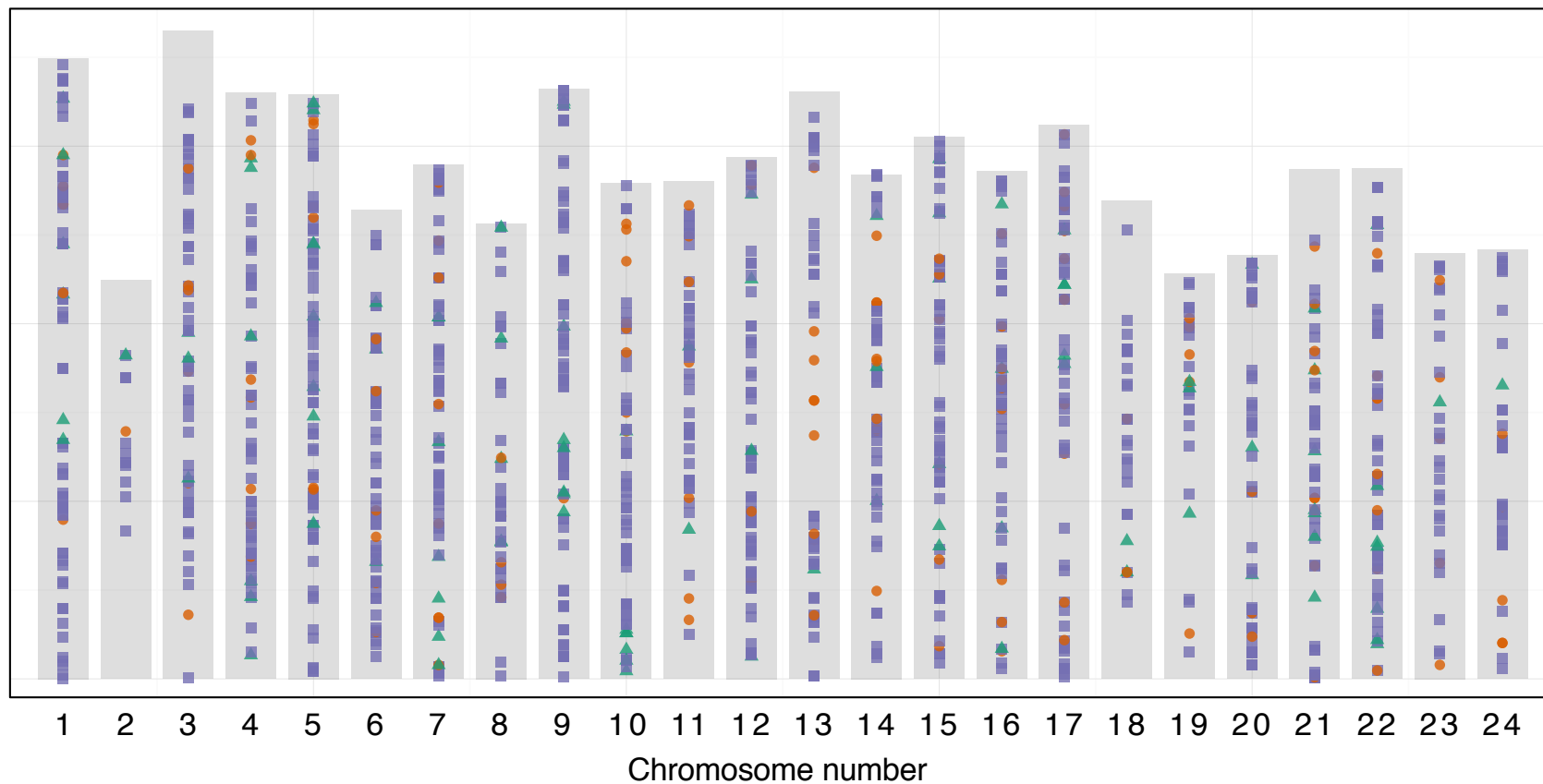


Figure 4



| | Mark | Chr1 | Chr2 | Chr3 | Chr4 | Chr5 | Chr6 | Chr7 | Chr8 | Chr9 | Chr10 | Chr11 | Chr12 | Total* ¹ |
|-------------------|------|------|------|------|------|------|------|------|------|------|-------|-------|-------|---------------------|
| NJPN-Tajima-Tango | ■ | 62 | 11 | 52 | 54 | 70 | 53 | 65 | 32 | 75 | 48 | 44 | 61 | 627 |
| SJPN-Tajima-Tango | ● | 7 | 1 | 6 | 7 | 6 | 6 | 8 | 5 | 1 | 8 | 7 | 4 | 66 |
| Tajima-Tango only | ▲ | 6 | 1 | 3 | 6 | 8 | 3 | 7 | 5 | 8 | 7 | 2 | 4 | 60 |

| | Mark | Chr13 | Chr14 | Chr15 | Chr16 | Chr17 | Chr18 | Chr19 | Chr20 | Chr21 | Chr22 | Chr23 | Chr24 | Total* ² |
|-------------------|------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|---------------------|
| NJPN-Tajima-Tango | ■ | 39 | 47 | 60 | 50 | 51 | 28 | 21 | 44 | 42 | 52 | 26 | 36 | 496 |
| SJPN-Tajima-Tango | ● | 10 | 9 | 6 | 9 | 13 | 3 | 5 | 5 | 9 | 8 | 5 | 6 | 88 |
| Tajima-Tango only | ▲ | 1 | 3 | 6 | 4 | 5 | 2 | 3 | 3 | 7 | 7 | 1 | 1 | 43 |

Figure 5

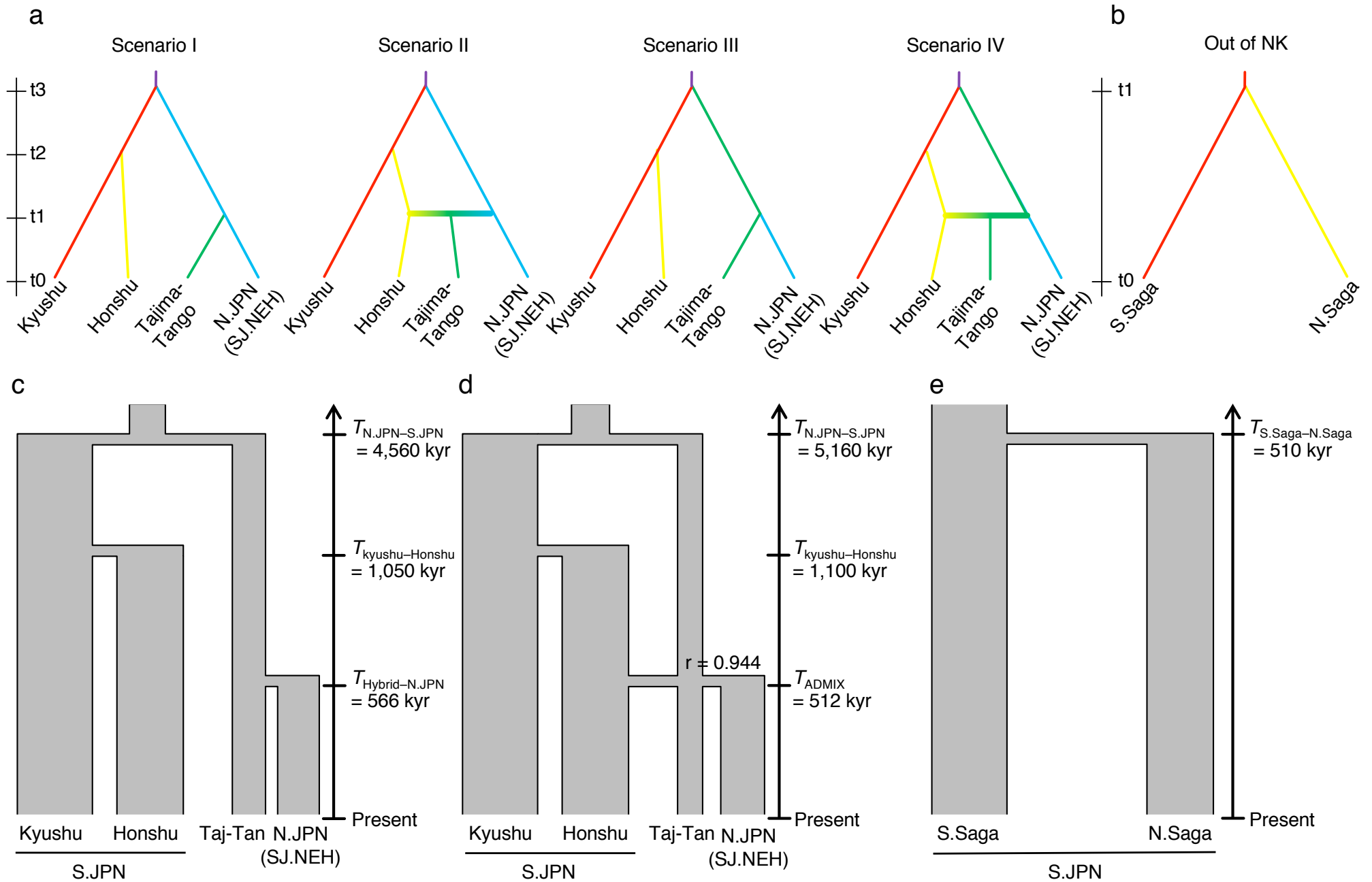


Figure 6

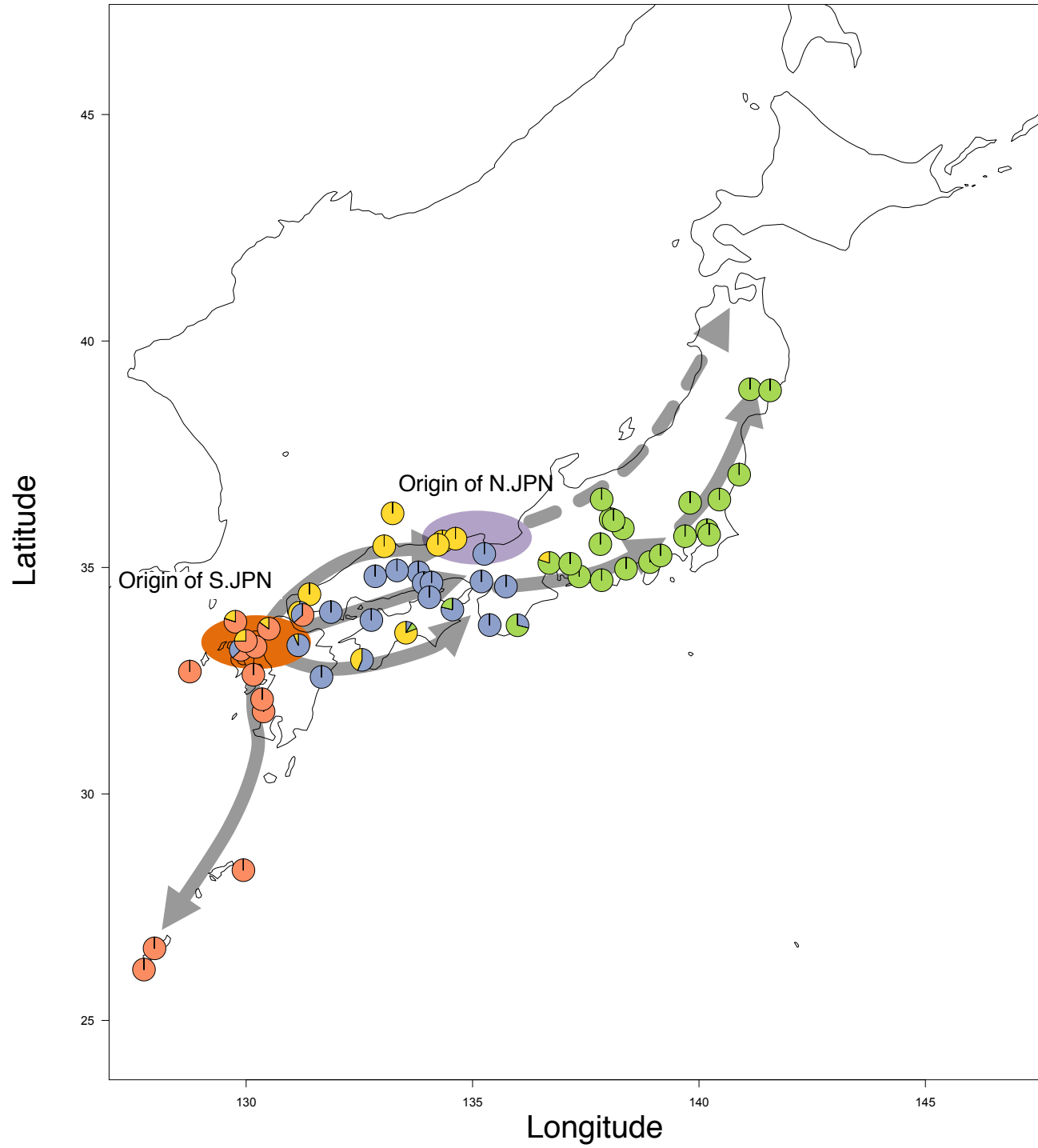


Table 1 Summary genetic statistics for five populations using "PopStat" dataset. These statistics include the mean number of individuals genotyped at each locus (*1), the number of variable sites in five populations (Variants), the number of polymorphic sites in each population (Polymorphic) and the number of variable sites unique to each population (Private). The number in parenthesis is Standard Error.

| Group | Number of Individual*1 | Total length of GBS loci (bp) | Variant | Polymorphic | Private | Major Allele Frequency | Observed Heterozygosity | Nucleotide Diversity |
|-----------------|------------------------|-------------------------------|---------|-------------|---------|------------------------|-------------------------|-----------------------|
| N.JPN (S.J.NEH) | 11.5 (+/-0.03) | 45968 | 2453 | 256 | 136 | 0.983 (+/-0.0014) | 0.023 (+/-0.0022) | 0.0014 (+/-0.0001) |
| Tajima -Tango | 4.4 (+/-0.01) | 44011 | 2453 | 199 | 71 | 0.980 (+/-0.0016) | 0.024 (+/-0.0022) | 0.0017 (+/-0.0001) |
| S.JPN | 41.6 (+/-0.08) | 45968 | 2453 | 1448 | 1231 | 0.958 (+/-0.0017) | 0.029 (+/-0.0016) | 0.0036 (+/-0.0001) |
| W.KOR | 2.4 (+/-0.01) | 44990 | 2453 | 467 | 342 | 0.940 (+/-0.0027) | 0.054 (+/-0.0032) | 0.0053 (+/-0.0003) |
| E.KOR | 2.4 (+/-0.01) | 44011 | 2453 | 346 | 203 | 0.953 (+/-0.0025) | 0.041 (+/-0.0030) | 0.0042 (+/-0.0002) |

Table 2 Demographic parameters estimated by DIYABC under three scenarios. Performance of each estimation was evaluated by RRMISE (the square Root of the Relative Mean Integrated Square Error), RMeanAD (the Relative Mean Absolute Deviation) and RRMSE (the square Root of the Relative Mean Square Error), which are output from DIYABC option "Compute bias and mean square error." The numbers with or without parentheses in the columns "Performances for estimating posterior distributions of parameters" are those of the statistics computed from the prior or posterior distribution of parameters, respectively.

| Scenario | Parameters | Type | Prior distribution | | | | Posterior parameter estimates | | | Performances for estimating posterior distributions of parameters | | |
|-----------|---------------------|---------|--------------------|------------|-----------|-----------|-------------------------------|-----------------------|----------------|---|-----------------|------------------|
| | | | min. | max. | mean | S.D. | Mean | 95% credible interval | | RRMISE | RMeanAD | RRMSE |
| II | N_{N_JPN} | Normal | 10,000 | 3,000,000 | 450,000 | 500,000 | 825,000 | 492,000 | 1,240,000 | 0.415 (-0.586) | 0.302 (-0.485) | 0.302 (-0.288) |
| | $N_{Tajima-Tango}$ | Normal | 10,000 | 3,000,000 | 600,000 | 500,000 | 613,000 | 369,000 | 901,000 | 0.35 (-0.808) | 0.256 (-0.61) | 0.251 (-0.358) |
| | N_{Kyushu} | Normal | 10,000 | 3,000,000 | 1,225,000 | 500,000 | 1,210,000 | 711,000 | 1,830,000 | 0.364 (-0.524) | 0.264 (-0.396) | 0.267 (-0.282) |
| | N_{honshu} | Normal | 10,000 | 3,000,000 | 1,050,000 | 500,000 | 1,490,000 | 974,000 | 2,020,000 | 0.291 (-0.442) | 0.219 (-0.368) | 0.213 (-0.289) |
| | $N_{Ancestor}$ | Uniform | 10,000 | 30,000,000 | - | - | 432,000 | 59,700 | 922,000 | 5.899 (-133.368) | 1.846 (-59.034) | 5.119 (-115.475) |
| | $T_{N_JPN-Hybrid}$ | Uniform | 10 | 4,000,000 | - | - | 566,000 | 311,000 | 908,000 | 0.409 (-2.417) | 0.300 (-1.725) | 0.266 (-1.589) |
| | $T_{kyushu-Honshu}$ | Uniform | 10 | 4,000,000 | - | - | 1,050,000 | 442,000 | 2,290,000 | 0.575 (-2.398) | 0.374 (-1.875) | 0.284 (-2.114) |
| | $T_{N_JPN-S_JPN}$ | Normal | 1,000,000 | 30,000,000 | 4,000,000 | 5,000,000 | 4,560,000 | 2,820,000 | 6,810,000 | 0.329 (-1.079) | 0.245 (-0.776) | 0.229 (-0.743) |
| III | N_{N_JPN} | Normal | 10,000 | 3,000,000 | 450,000 | 500,000 | 812,000 | 390,000 | 1,330,000 | 0.447 (-0.635) | 0.332 (-0.506) | 0.303 (-0.326) |
| | $N_{Tajima-Tango}$ | Normal | 10,000 | 3,000,000 | 600,000 | 500,000 | 480,000 | 266,000 | 739,000 | 0.375 (-1.208) | 0.284 (-0.88) | 0.261 (-0.736) |
| | N_{Kyushu} | Normal | 10,000 | 3,000,000 | 1,225,000 | 500,000 | 1,290,000 | 758,000 | 1,910,000 | 0.356 (-0.483) | 0.265 (-0.373) | 0.265 (-0.252) |
| | N_{honshu} | Normal | 10,000 | 3,000,000 | 1,050,000 | 500,000 | 1,450,000 | 913,000 | 2,020,000 | 0.316 (-0.446) | 0.235 (-0.369) | 0.235 (-0.279) |
| | $N_{Ancestor}$ | Uniform | 10,000 | 30,000,000 | - | - | 390,000 | 57,200 | 849,000 | 6.947 (-131.19) | 2.285 (-62.259) | 5.997 (-113.55) |
| | T_{ADMIX} | Uniform | 10 | 4,000,000 | - | - | 512,000 | 242,000 | 842,000 | 0.397 (-4.374) | 0.308 (-2.163) | 0.249 (-3.309) |
| | $T_{kyushu-Honshu}$ | Uniform | 10 | 4,000,000 | - | - | 1,100,000 | 446,000 | 2,420,000 | 0.576 (-2.256) | 0.381 (-1.76) | 0.296 (-1.98) |
| | $T_{N_JPN-S_JPN}$ | Normal | 1,000,000 | 30,000,000 | 4,000,000 | 5,000,000 | 5,160,000 | 3,080,000 | 7,780,000 | 0.345 (-0.916) | 0.262 (-0.65) | 0.239 (-0.588) |
| r | Uniform | 0.001 | 0.999 | - | - | 0.944 | 0.897 | 0.987 | 0.027 (-0.559) | 0.02 (-0.472) | 0.014 (-0.469) | |
| Out of NK | $N_{S.Saga}$ | Normal | 10,000 | 3,000,000 | 1,300,000 | 500,000 | 1,326,000 | 960,075 | 1,670,250 | 0.178 (-0.31) | 0.138 (-0.241) | 0.120 (-0.095) |
| | $N_{N.Saga}$ | Normal | 10,000 | 3,000,000 | 1,250,000 | 500,000 | 1,224,000 | 892,500 | 1,593,750 | 0.181 (-0.322) | 0.142 (-0.251) | 0.119 (-0.096) |
| | $T_{S.Saga-N.Saga}$ | Uniform | 1,000 | 4,000,000 | - | - | 510,000 | 337,875 | 679,575 | 0.224 (-5.918) | 0.175 (-4.488) | 0.143 (-4.429) |