

# 1 Oceans apart: Heterogeneous patterns of parallel 2 evolution in sticklebacks

3 Bohao Fang\*, Petri Kemppainen\*, Paolo Momigliano, Xueyun Feng and Juha Merilä

4  
5 *Ecological Genetics Research Unit, Organismal and Evolutionary Biology Research*  
6 *Program, Faculty of Biological and Environmental Sciences, University of Helsinki, FI-*  
7 *00014 Helsinki, Finland*

8 \*These authors contributed equally to this work

9 Correspondence to: PK and BF

10 E-mail: [petri.kemppainen@helsinki.fi](mailto:petri.kemppainen@helsinki.fi); [bohao.fang@helsinki.fi](mailto:bohao.fang@helsinki.fi)

11

## 12 **Abstract**

13 An important model system for the study of genomic mechanisms underlying parallel  
14 ecological adaptation in the wild is the three-spined stickleback (*Gasterosteus aculeatus*),  
15 which has repeatedly colonized and adapted to freshwater from the sea throughout the  
16 northern hemisphere. Previous studies have identified numerous genomic regions showing  
17 consistent genetic differentiation between freshwater and marine ecotypes, but these are  
18 typically based on limited geographic sampling and are biased towards studies in the  
19 Eastern Pacific. We analysed population genomic data from marine and freshwater  
20 ecotypes of three-spined sticklebacks with from a comprehensive global collection of  
21 marine and freshwater ecotypes to detect loci involved in parallel evolution at different  
22 geographic scales. Our findings highlight that most signatures of parallel evolution were

23 unique to the Eastern Pacific. Trans-oceanic marine and freshwater differentiation was  
24 only found in a very limited number of genomic regions, including three chromosomal  
25 inversions. Using both simulations and empirical data, we demonstrate that this is likely  
26 due to both the stochastic loss of freshwater-adapted alleles during founder events during  
27 the invasion of the Atlantic basin and selection against freshwater-adapted variants in the  
28 sea, both of which have reduced the amount of standing genetic variation available for  
29 freshwater adaptation outside the Eastern Pacific region. Moreover, the existence of highly  
30 elevated linkage disequilibrium associated with marine-freshwater differentiation in the  
31 Eastern Pacific is also consistent with a secondary contact scenario between marine and  
32 freshwater populations that have evolved in isolation from each other during past glacial  
33 periods. Thus, contrary to what earlier studies focused on Eastern Pacific populations  
34 have led us to believe, parallel marine-freshwater differentiation in sticklebacks is far less  
35 prevalent and pronounced in all other parts of the species global distribution range.

36

37 **Keywords:** *Gasterosteus aculeatus*; genetic differentiation; linkage disequilibrium; local  
38 adaptation; parallel evolution

## 39 Introduction

40

41 The extent to which the evolution of similar phenotypes arises by selection acting on  
42 shared ancestral polymorphism (i.e. parallel evolution (Schluter & Conte 2009)) or via  
43 distinct molecular evolutionary pathways (i.e. convergent evolution (Arendt & Reznick  
44 2008; DeFaveri *et al.* 2011; Stern 2013) is a major question in evolutionary biology. A  
45 powerful approach to disentangle these processes is to study the genomic architecture  
46 underlying the repeated evolution of similar phenotypes in similar environments. After the  
47 retreat of Pleistocene glaciers, marine three-spined sticklebacks (*Gasterosteus aculeatus*)  
48 colonized and adapted to many newly formed freshwater habitats by adopting similar  
49 changes in a number of morphological, physiological, life history and behavioural traits  
50 (Bell & Foster 1994; Gibson 2005; Hendry *et al.* 2013; Lescak *et al.* 2015; Östlund-Nilsson  
51 *et al.* 2006). Thus, this species has become one of the most widely used model systems to  
52 study the molecular basis of adaptive evolution in vertebrates in the wild (McKinnon &  
53 Rundle 2002).

54 Previous studies of the three-spined stickleback model system have quantified the extent  
55 of parallel evolution by identifying genomic regions that are consistently differentiated  
56 between marine and freshwater ecotypes sampled across different geographic areas  
57 (DeFaveri *et al.* 2011; Ferchaud & Hansen 2016; Hohenlohe *et al.* 2010; Hohenlohe &  
58 Magalhaes 2019; Jones *et al.* 2012; Liu *et al.* 2018; Pujolar *et al.* 2017; Terekhanova *et al.*  
59 2019; Terekhanova *et al.* 2014). The focus has historically been on the Eastern Pacific  
60 region (Chan *et al.* 2010; Colosimo *et al.* 2005; Hohenlohe *et al.* 2010; Hohenlohe &  
61 Magalhaes 2019; Jones *et al.* 2012; Nelson & Cresko 2018), but several recent studies  
62 have focused on Atlantic populations (Ferchaud & Hansen 2016; Liu *et al.* 2018; Pujolar *et*  
63 *al.* 2017; Terekhanova *et al.* 2019; Terekhanova *et al.* 2014). However, only two studies

64 have thus far included samples from a larger (global) geographic range (DeFaveri *et al.*  
65 2011; Jones *et al.* 2012). Based on whole genome sequence data from a limited number  
66 of individuals from Eastern Pacific and Atlantic populations (n = 21), Jones *et al.* (2012)  
67 identified ~200 genomic regions that consistently separated marine and freshwater  
68 individuals globally, representing roughly 0.5% of the dataset. They also found that 2.83%  
69 of the genome showed signatures of parallel selection in Eastern Pacific freshwater  
70 locations – approximately six times more than that at the global scale (tree i,  
71 Supplementary Fig. 2 and Supplementary Table 7 in Jones *et al.* (2012) – suggesting that  
72 more loci contribute to parallel evolution at smaller geographic (regional) scales. However,  
73 since this pattern was unique to the Eastern Pacific (and the focus was on global  
74 parallelism) its implications for a holistic understanding of marine-freshwater differentiation  
75 at both regional and global scales was never discussed. Such global heterogeneous  
76 ecotype divergence is consistent with the results of several other studies as well. Focusing  
77 on 26 candidate genes in six pairs of marine-freshwater populations across the globe,  
78 DeFaveri *et al.* (2011) found that only ~50% of the genes under divergent selection were  
79 shared across more than three population pairs, and none were shared among all  
80 populations. This suggested a limited re-reuse of ancestral polymorphism at the global  
81 scale, implicating either an important role of convergent evolution at larger geographic  
82 scales (DeFaveri *et al.* 2011), or geographic heterogeneity in selective pressure among  
83 different (DeFaveri *et al.* 2011) freshwater ecosystems (DeFaveri *et al.* 2011; Stern 2013).  
84 Furthermore, studies focusing on parallel evolution within oceans, and even smaller  
85 geographic regions, show striking differences in the proportion of loci involved in parallel  
86 freshwater adaptation between Pacific and Atlantic regions (Ferchaud & Hansen 2016;  
87 Hohenlohe *et al.* 2010; Jones *et al.* 2012; Liu *et al.* 2018; Nelson & Cresko 2018; Pujolar *et*  
88 *al.* 2017; Terekhanova *et al.* 2019; Terekhanova *et al.* 2014). For instance, Terekhanova *et*

89 *al.* (2014, 2019) recovered only 21 highly localized genomic regions involved in parallel  
90 freshwater differentiation in the White Sea, in contrast to Hohenlohe *et al.* (2010) and  
91 Nelson & Cresko (2018) who found large genomic regions involved in parallel freshwater  
92 differentiation across almost all chromosomes in the Eastern Pacific populations.  
93 Therefore, the potential mechanisms underlying this apparent large-scale geographic  
94 heterogeneity in genome-wide patterns of parallel evolution in three-spined sticklebacks  
95 remain unexplored. To this end, we analysed population genomic data from a  
96 comprehensive sampling of all major geographic areas inhabited by the three-spined  
97 stickleback, and employed unsupervised and supervised methods to detect loci involved in  
98 parallel marine-freshwater differentiation at different geographical scales. Based on earlier  
99 observations (DeFaveri *et al.* 2011; Ferchaud & Hansen 2016; Kempainen *et al.* 2015;  
100 Liu *et al.* 2018; Pujolar *et al.* 2017; Terekhanova *et al.* 2019; Terekhanova *et al.* 2014), we  
101 hypothesize that the genetic parallelism in response to freshwater colonization by marine  
102 sticklebacks is heterogeneous at the global scale, and that the degree of genetic  
103 parallelism is much stronger in the Eastern Pacific region than elsewhere.

104 We further seek to understand and discuss the ultimate causes of the marked regional  
105 differences in genome-wide signatures of parallel genetic differentiation among ecotypes.  
106 To explain the mechanism behind the repeated use of the same alleles in independent  
107 freshwater populations of sticklebacks, Schluter & Conte (2009) proposed the “transporter  
108 hypothesis”. This hypothesis postulates that three-spined sticklebacks have repeatedly  
109 colonized and adapted to freshwater environments via selection on standing genetic  
110 variation in large marine populations. These freshwater-adapted alleles are in turn  
111 maintained in the marine populations by recurrent gene flow with previously colonized  
112 freshwater populations. Three-spined stickleback populations have persisted in the  
113 Eastern Pacific for approximately 26 Mya (Betancur *et al.* 2015; Matschiner *et al.* 2011;

114 Meynard *et al.* 2012; Sanciangco *et al.* 2016) and from there recolonized the Western  
115 Pacific and Atlantic Ocean basins following local extinctions much more recently, during  
116 the late Pleistocene (36.9-346.5 kya; Fang *et al.* 2020; Fang *et al.* 2018; Orti *et al.* 1994).  
117 During bottlenecks and founder events, rare alleles are lost at a higher rate than common  
118 alleles (Halliburton & Halliburton 2004; Hyten *et al.* 2006). Since freshwater-adapted  
119 alleles exist in the marine populations only at low frequencies (Schluter & Conte 2009), it is  
120 likely that they were lost to a higher degree than neutral variation during geographic range  
121 expansions from the Eastern Pacific (via the sea), thereby reducing the amount of  
122 standing genetic variation available for freshwater adaptation outside of the Eastern  
123 Pacific. To test this hypothesis, we used individual-based forward simulations designed to  
124 mimic the transporter hypothesis, and the general global population demographic history  
125 of three-spined sticklebacks outlined above. We conclude with a discussion on other  
126 potential biological and demographic explanations for the high degree of geographic  
127 heterogeneity in patterns of parallel genomic differentiation, and reflect upon the  
128 representativeness of the Eastern Pacific three-spined stickleback populations as a  
129 general model for the study of parallel evolution.

130

131

## 132 **Material and Methods**

### 133 **Sample collection**

134 We obtained population genomic data from 166 individuals representing both marine and  
135 freshwater ecotypes from the Eastern and Western Pacific, as well as from the Eastern  
136 and Western Atlantic Oceans (Fig. 1i, Supplementary Table 1 and Supplementary Fig. 1).  
137 Additional data from previously published studies were retrieved from GenBank. Fish

138 collected for this study were sampled with seine nets, minnow traps and electrofishing.  
139 Specimens were preserved in ethanol after being euthanized with an overdose of Tricaine  
140 mesylate (MS222).

141 To study the extent of genetic parallelism among freshwater sticklebacks with different  
142 phylogeographic histories, we classified global samples into seven biogeographic regions  
143 based on their phylogenetic affinities: (i) Eastern Pacific, (ii) Western Pacific, (iii) Western  
144 Atlantic, (iv) White and Barents Seas, (v) North Sea and British Isles, (vi) Baltic Sea and  
145 (vii) Norwegian Sea (Fang et al., 2018; Fig. 1i). A summary of coordinates, ecotype and  
146 population information on the sampled individuals and re-acquired samples is given in the  
147 Supplementary Table 1.

148

#### 149 **Sequencing and genotype likelihood estimation**

150 Restriction site associated DNA sequencing (RADseq) using the enzyme *PstI* was  
151 performed for the 62 individuals sampled in this study, using the same protocol as in Fang  
152 et al. (2018), where DNA library preparation and sequencing method are described in  
153 detail. The raw RAD sequencing data has been uploaded to GenBank (Accessions  
154 SAMN14078677-SAMN14078738). Previously published RADseq and whole genome  
155 sequencing (WGS) data for an additional 104 individuals from 62 populations were  
156 retrieved from GenBank (Supplementary Table 1). All RADseq and WGS datasets were  
157 mapped to the three-spined stickleback reference genome (release-92, retrieved from  
158 Ensembl (Hubbard *et al.* 2005) using BWA mem v0.7.17 (Li & Durbin 2009). PCR  
159 duplicates were removed using the program Stacks v2.5 (Catchen *et al.* 2013) for pair-end  
160 RAD data, and SAMtools v1.9 (function “markdup”(Li 2011)) for whole genome data. Given  
161 the heterogeneity in sequencing depth among different datasets, and particularly the very

162 low coverage of the data retrieved from Jones *et al.* (Jones *et al.* 2012), most of the  
163 analyses were performed directly using genotype likelihoods, avoiding variant calling  
164 whenever possible. Genotype likelihoods were estimated from the mapped reads using  
165 the model of SAMtools (Li 2011) as implemented in the program suite ANGSD v0.929  
166 (Korneliussen *et al.* 2014). Full scripts for the genotyping and filtering parameters are  
167 publically available through DRYAD (doi: xxx). Bases with a q-score below 20 (-minQ 20)  
168 and reads with mapping quality below 25 (-minMapQ 25) were removed, and variants were  
169 only retained if they had a p-value smaller than  $1e^{-6}$  (-SNP\_pval  $1e^{-6}$  flag in ANGSD). We  
170 retained sites with a minimum read depth of two (-minIndDepth 2) in at least 80% of the  
171 sampled individuals (-minInd 133). The sex chromosome (Chr. XIX (Kitano *et al.* 2009;  
172 Natri *et al.* 2013)) was excluded from downstream analyses due to sex-specific genomic  
173 heterogeneity (Hedrick 2007; Schaffner 2004). The raw output of genotype likelihoods  
174 from all 166 individuals comprised 2,511,922 genome-wide loci.

175

## 176 **Unsupervised approach to determine marine-freshwater differentiation**

177 We conducted Linkage Disequilibrium Network Analysis (LDna) on the whole dataset  
178 (2,511,922 SNPs) to identify and extract clusters of highly correlated loci, i.e. sets of loci  
179 affected by the same evolutionary processes. LDna uses a pairwise matrix of LD values,  
180 estimated by  $r^2$ , to produce a single linkage clustering tree. The hierarchical clustering  
181 algorithm uses the LD matrix to combine two clusters connected to each other by at least  
182 one edge. In the resulting tree, the nodes represent clusters of loci connected by LD  
183 values above thresholds, where the threshold value is proportional to the distance from the  
184 root (Kempainen *et al.* 2015). As the LD threshold is sequentially lowered, an increasing  
185 number of loci will be connected to each other in a fashion that reflects their similarity in  
186 phylogenetic signals. For each cluster merger (with decreasing LD threshold), the change



187 in median LD between all pairwise loci in a cluster before and after the merger is estimated  
188 as  $\lambda$ . When two highly interconnected clusters merge,  $\lambda$  will be large (unlike when only a  
189 single locus is added to an existing cluster), signifying that these two clusters bear distinct  
190 phylogenetic signals. LDna is currently limited to ~20,000 SNPs at a time due to its  
191 dependence on LD estimates for all pairwise comparisons between loci in the dataset. To  
192 analyse the whole dataset, we applied a novel three-step LDna approach to reduce the  
193 complexity of the data in a nested fashion. First, we started with non-overlapping windows  
194 within each chromosome (Li *et al.* 2018), then performed the analysis on each  
195 chromosome individually, and finally on the whole dataset (Supplementary Information 1).  
196 In all steps of LDna, we estimated LD between loci from genotype likelihoods using the  
197 program ngsLD (Fox *et al.* 2019), setting the minimum SNP minor allele frequency at 0.05.  
198 Full scripts for the LDna analyses are provided in DRYAD (doi: xxx). In the first step, we  
199 only kept loci that were in high LD with at least one locus ( $r_2 > 0.8$ ) within a window of 100  
200 SNPs, as most SNPs in the data were not correlated with any other adjacent loci (so-  
201 called singleton clusters (Li *et al.* 2018), and thus, are unlikely to be informative in the  
202 LDna analyses.

203 The main evolutionary phenomena that cause elevated LD between large sets of loci in  
204 population genomic datasets are polymorphic inversions, population structure and local  
205 adaptation, all of which are expected to be present in our dataset (Kemppainen *et al.*  
206 2015). There are specific and distinct predictions about the population genetic signal and  
207 the distribution of loci in the genome that arise from these evolutionary phenomena  
208 (Kemppainen *et al.* 2015). First, clusters with LD signals caused by inversions are  
209 expected to predominantly map to the specific genomic region where the inversion is  
210 situated. In addition, a Principal Component Analysis (PCA) of these loci is expected to  
211 separate individuals based on karyotype. In general, the heterokaryotype would be

212 intermediate to the two alternative homokaryotypes (provided that all karyotypes exist in  
213 the dataset), and the heterokaryotypic individuals would show higher observed  
214 heterozygosity than the homokaryotypes. However, this is not always so clear, for instance  
215 when the inversion is new and mutational differences have not yet accumulated. Second,  
216 a PCA based on loci whose frequencies are shaped by genetic drift is expected to  
217 separate individuals on the basis of geographic location, with no (or very little) separation  
218 between marine and freshwater ecotypes. Third, an LD signal caused by local adaptation  
219 (globally) is expected to cluster individuals based on ecotype, regardless of geographic  
220 location, with both the locus distribution and LD patterns being negatively correlated with  
221 local recombination rate to some extent (Roesti *et al.* 2014; Roesti *et al.* 2013). The reason  
222 for this correlation is that gene flow between ecotypes erodes genetic differentiation in  
223 sites linked to locally adapted loci with the exception of regions where recombination is  
224 restricted (for instance in inversions, or close to centromeres or telomeres). No such  
225 pattern is expected for LD caused by population structuring, as the main source of this LD  
226 is the random genetic drift that, in the absence of gene flow, generates LD in a fashion that  
227 is independent of genome position (Kempainen *et al.* 2015) and background selection is  
228 not expected to result in strong patterns of genetic differentiation across genomes  
229 (Matthey-Doret & Whitlock 2019; Stankowski *et al.* 2019). If a set of loci contributes to local  
230 adaptation exclusively in a particular geographic area, a PCA based on these loci will only  
231 separate individuals based on ecotype in that region. We considered loci to be involved in  
232 parallel evolution only if they grouped individuals of the same ecotype from more than one  
233 independent location. Otherwise, it is not possible to discern drift from local adaptation,  
234 particularly if  $N_e$  is small (i.e. genetic drift is strong). To determine if an LD-cluster was  
235 likely associated with parallel freshwater differentiation, we first used expectation  
236 maximisation (EM) and hierarchical clustering methods to identify clusters of individuals in

237 PCAs that contained a minimum of seven individuals, of which at least 90% are freshwater  
238 ecotypes (the “in-group”; dotted line; Fig. 1a-h, Supplementary Fig. 2). Second, if such in-  
239 groups were detected, we used permutations to further test whether this cluster contained  
240 more freshwater individuals than expected by chance (Supplementary Information 2). With  
241 less than seven in-group individuals, there was no power to detect significant associations,  
242 even if all individuals were freshwater ecotypes. We benchmarked LDna by quantifying the  
243 proportion of regions previously identified by Jones et al. (2012) as involved in marine-  
244 freshwater ecotype differentiation (globally and within the Eastern Pacific) that were  
245 correctly recovered by LD-cluster loci (Supplementary Information 5).

246

#### 247 **Supervised approaches to determine marine-freshwater differentiation**

248 Genome-wide allelic differentiation ( $F_{ST}$  estimated from genotype likelihoods in ANGSD)  
249 between marine and freshwater ecotypes was estimated separately for the three major  
250 oceans in our study: Eastern Pacific, Western Pacific and Atlantic Oceans. All available  
251 samples were always used, but due to the small number of available Pacific marine  
252 individuals, marine individuals from the Eastern ( $n=4$ ) and Western Pacific ( $n=13$ ) were  
253 pooled and treated as a combined Pacific marine group in Eastern and Western Pacific  
254 ecotype comparisons (Supplementary Table 2) in order to reduce the noise of the SNP-  
255 based analyses. To determine whether pooling Eastern and Western Pacific marine  
256 individuals could bias  $F_{ST}$  estimates, we first estimated  $F_{ST}$  in 100 kb windows for  
257 Eastern Pacific marine vs. Eastern Pacific freshwater and Western Pacific marine vs.  
258 Eastern Pacific freshwater individuals, as using large windows allowed us to get precise  
259 estimates even if we sampled only four marine individuals from the Eastern Pacific. The  
260 two sets of window-based pairwise  $F_{ST}$  estimates were highly correlated ( $r = 0.904$ ;  
261  $P < 0.001$ ; Supplementary Fig. 4d-g), suggesting that pooling marine individuals from

262 Eastern and Western Pacific should not strongly affect SNP based estimates. Note that  
263 from the results of the unsupervised LDna, two Eastern Pacific freshwater individuals from  
264 Kodiak Island, Alaska (ALA population) never grouped with the other Eastern Pacific  
265 freshwater individuals. Therefore, in agreement with earlier phylogenetic analyses (Fang *et*  
266 *al.* 2018), these two individuals were excluded from the supervised analyses. The squared  
267 correlation coefficient of  $F_{ST}$  before and after this exclusion was 0.88, indicating that this  
268 did not affect the results.

269 In each comparison, the sites were firstly filtered from raw mapped reads, retaining sites  
270 with less than 25% missing data with quality control (-minIndDepth 1, -uniqueOnly 1, -  
271 remove\_bads 1, -minMapQ 20, -minQ 20). We retained only variable sites (-SNP\_pval 1e-  
272 6) in each region, resulting in 1,218,858 SNPs in the Eastern Pacific, 1,072,257 SNPs in  
273 the Western Pacific, and 1,681,923 SNPs in the Atlantic Ocean. We then obtained  
274 genotype likelihoods and site allele frequency likelihoods of the variants (-GL 1, -doSaf 1).  
275 Based on these likelihoods, we estimated the two-dimensional site-frequency spectrum  
276 (SFS) for each pair of ecotypes (realSFS) and calculated the pairwise weighted  $F_{ST}$   
277 (realSFS fst).

278

## 279 **Proof of concept using simulated data**

280 Several potential explanations for geographic heterogeneity in parallel patterns of marine-  
281 freshwater differentiation in three-spined sticklebacks have been suggested (DeFaveri *et*  
282 *al.* 2011). One such explanation that has not received much attention in the context of  
283 three-spined sticklebacks is the stochastic loss of freshwater-adapted alleles due to  
284 founder events when three-spined sticklebacks colonized the rest of the world from the  
285 Eastern Pacific in the late Pleistocene (see Introduction). Thus, as a proof of concept, we

286 used forward-in-time simulations performed with quantiNemo (Neuenschwander *et al.*  
287 2008) to investigate the conditions under which parallel islands of differentiation between  
288 marine and freshwater ecotypes can arise under such a scenario.

289 Our simulations were aimed at recreating the transporter hypothesis model in the Eastern  
290 Pacific (referred to as “*Pac*” in the context of simulations), to simulate the colonization of  
291 the Atlantic (referred to as “*Atl*” in the context of simulations) from *Pac* 60-30 Kya during  
292 the last known opening of the Bering Strait (Fang *et al.* 2020; Hu *et al.* 2010; Meiri *et al.*  
293 2014) and the subsequent post-glacial (10 Kya) colonization of newly formed freshwater  
294 habitats in both oceans (simulation details can be found in Supplementary Information 4).  
295 In short, simulations begin with one marine population in *Pac* connected to five  
296 independent freshwater populations by symmetrical gene flow (i.e. no gene flow exists  
297 between any of the freshwater populations; Fig. 3a) for 10k generations (40-50 kya). This  
298 is followed by colonisation of *Atl* from *Pac* (with *Atl* having identical population structure to  
299 *Pac*) by allowing one or five migrants per generation between the oceans for 2 Ky (38-40  
300 Kya), after which no further gene flow is possible. The retreat of the Pleistocene  
301 continental ice sheets (at 10 Kya; Fig. 3d) and the colonization of newly formed freshwater  
302 habitats is simulated by removing four of the freshwater populations, immediately followed  
303 by the emergence of four new (post-glacial) freshwater populations (in both *Pac* and *Atl*;  
304 Fig. 3e). The fifth freshwater population remains as a “glacial refugia” that continues to  
305 feed freshwater-adapted alleles to the sea as standing genetic variation. Post-glacial local  
306 adaptation is thus only possible due to the spread of freshwater-adapted alleles from the  
307 sea in accordance with the transporter hypothesis (Schluter & Conte, 2009; Fig. 3a-e).

308 Marine-freshwater differentiation was based on bi-allelic QTL with allelic effects of either  
309 zero or ten, with the selection optima in the marine habitat being zero and the selection  
310 optima in all freshwater populations being 20. Thus, a freshwater individual homozygous

311 for allele 2 for a given QTL meant that the individual was at its optimal phenotype, and *vice*  
312 *versa* for marine individuals. Selection intensities were such that a sufficient amount of  
313 standing genetic variation was allowed in the sea and rapid local adaptation in freshwater  
314 was possible (see Supporting Information 4 for details). In simulations, all allele  
315 frequencies started from 0.5 in all populations (including the QTL in the freshwater  
316 habitats). The simulated genome was comprised of ten equally sized chromosomes, with a  
317 total genome size of 1000 bps. Regions of both low (centromeric regions) and high  
318 (chromosome arms) recombination were represented (Supplementary Information 4).  
319 Either 3, 6 or 9 QTL per chromosome were randomly placed in eight of the chromosomes,  
320 after which the positions were fixed, leaving the last two chromosomes without any QTL.  
321 Twenty replicate simulations were run for each of the six different parameter settings (two  
322 levels of trans-oceanic gene flow rates and three different QTL densities). The frequency  
323 of freshwater-adapted alleles was recorded at 50-generation intervals throughout the  
324 simulations. Population genomic data were saved at the end of the simulations  
325 (representing present-day sampling).

326

### 327 **Linking empirical data to simulated data**

328 LDna identified one major cluster (LD-cluster 2; see Results) that separated all Eastern  
329 Pacific freshwater individuals from the remaining individuals (Atlantic, Western Pacific and  
330 marine samples from the Eastern Pacific, pooled). From the simulated data, we first sub-  
331 sampled individuals from *Pac* and *Atl* to match the samples size of the empirical data  
332 (excluding the Western Pacific samples, as this ocean was not included in the  
333 simulations), and used LDna to detect clusters similar to LD-cluster 2 using Cluster  
334 Separation Scores (CSS; custom R-scripts available from DRYAD). Cluster Separation  
335 Scores were calculated as the Euclidean centroid distance in a PCA (based on

336 coordinates from the two first principal components scaled by their eigenvalues) between  
337 two groups of individuals, standardized by the longest distance between any two  
338 individuals in the PCA (CSS thus ranging between [0,1]). PCA of the simulated datasets  
339 were performed by the function `snpGdsPCA` from the R-package `SNPRelate` (Zheng *et al.*  
340 2012). CSS scores are known to correlate with  $F_{ST}$ , but give higher resolution when  
341 genetic differentiation is high, and are less sensitive to small sample sizes (Jones *et al.*  
342 2012). In LDna, we are only interested in clusters with high  $\lambda$ -values (see above and  
343 Supplementary Information 1). Therefore, from the ten LD-clusters with the highest  $\lambda$ -  
344 values (from each simulated dataset), we considered the cluster with the highest CSS  
345 between *PF* (Eastern Pacific freshwater) and non-*PF* individuals to be the strongest  
346 candidates for showing high ecotype differentiation specifically in *Pac*, and thus, the most  
347 similar to LD-cluster 2 in the empirical data. The non-*PF* individuals were comprised of *PM*  
348 (Eastern Pacific marine), *AF* (Atlantic freshwater) and *AM* (Atlantic marine) individuals  
349 pooled. To further assess how similar the patterns of population differentiation (in PCAs)  
350 were in the above LD-clusters (simulated data) and empirically obtained LD-cluster 2, we  
351 compared the CSS's for all pairwise comparisons between *PF* and the other three groups  
352 of individuals (i.e. "*PF vs. PM*", "*PF vs. AF*" and "*PF vs. AM*") in the simulated and  
353 empirical datasets. To further assess the extent to which clusters similar to LD-cluster 2  
354 could be produced in the *Atl*, we used the same procedure as above to look for the LD-  
355 clusters with the highest CSS scores between *AF* and non-*AF* individuals (*AM*, *PF* and  
356 *PM*), and calculated CSS scores between *AF* individuals and the three non-*AF* groups.

357

## 358 Results

359 **Marine-freshwater divergence determined by unsupervised and supervised**  
360 **approaches**

361 The first step of LDna on the empirical dataset (2,511,922 SNPs derived from 166  
362 individuals worldwide) identified 214,326 loci that were in high LD with at least one other  
363 locus within windows of 100 SNPs (Supplementary Information 1). The next step of  
364 performing LDna on each chromosome separately (only using one locus from each LD-  
365 cluster from step one; Supplementary information 1) resulted in 81 distinct LD-clusters.  
366 From these, a final 29 LD-clusters were obtained (pooling within chromosome LD-clusters  
367 whenever they were grouped by LDna in the final step; Supplementary information 1),  
368 containing a total of 71,064 loci (*viz.* Cluster 1-29). Eight of these LD-clusters associated  
369 with geographic structure and genetic parallelism are highlighted in Fig. 2a-h. Details of all  
370 29 clusters can be found in Supplementary Fig. 2 and Supplementary Table 3.

371 LDna successfully recovered most of the previously identified regions from Jones et al.  
372 (2012) that differentiated marine from freshwater ecotypes, and failed to recover only small  
373 regions that had low coverage and relatively low levels of marine-freshwater differentiation  
374 (Supplementary Information 5 and Supplementary Fig. 3).

375

376 *Trans-oceanic marine-freshwater parallelism*

377 LD-clusters 5, 6, 10, 11, 12, 13, 16, 18, 20, 22, 25 and 27 (a total of 2,502 loci, 0.100% of  
378 the dataset; see four representatives in Fig. 1e-h, and Supplementary Fig. 2, Table 3)  
379 grouped multiple freshwater individuals from different geographic regions across the  
380 Pacific and Atlantic Oceans (for all  $P < 0.05$ , permutation test for ecotype differentiation),  
381 reflecting genetic marine-freshwater parallelism on a global (trans-oceanic) scale. Within  
382 those, LD-clusters 6, 11, 12, 16 and 27 (a total of 1,639 loci, 0.065% of the dataset)



383 similarly showed high marine-freshwater differentiation ( $F_{ST}$ , Fig 1e-h, Fig. 2c-e) in both the  
384 Eastern Pacific and Atlantic, further suggesting global parallelism. Particularly, loci from  
385 LD-cluster 11 mapped to four distinct regions on Chr. V of which one (Fig. 2c-e) mark the  
386 position of the *Ectodysplasin* (EDA) gene that is famously responsible for marine-  
387 freshwater differences in lateral armour plate development worldwide (Colosimo *et al.*  
388 2005). In contrast, although the remaining clusters (5, 10, 13, 18, 20, 22, 25, a total of 863  
389 loci, 0.034% of the dataset) also grouped freshwater individuals from both the Pacific and  
390 Atlantic Oceans (similar to LD-cluster 29 above; Fig. 2c-e and Supplementary Fig. 2), they  
391 showed much less marine-freshwater differentiation in the Atlantic Ocean than in the  
392 Eastern Pacific (Fig. 2c-e). Among the LD-clusters associated with marine-freshwater  
393 differentiation, LD-clusters 6 and 22 covered previously known chromosomal inversions on  
394 Chr. I and XI, respectively (Jones *et al.*, 2012, Fig. 1e, Supplementary Fig. 2 and  
395 Supplementary Table 3). In addition, we also found a putative novel inversion on Chr. V  
396 (LD-cluster 19,241 loci) that was not associated with marine-freshwater differentiation  
397 (Supplementary Fig. 2 and Table 3). Neither LDna nor  $F_{ST}$  analyses could detect any  
398 regions of marine-freshwater differentiation in the Western Pacific (Supplementary Fig. 4)  
399 and thus, this region is not considered further here.

400

#### 401 *Eastern Pacific marine-freshwater parallelism*

402 LD-clusters 2 (53,785 loci, 2.141% of the dataset, Fig. 1b) and 21 (183 loci, 0.007% of the  
403 dataset, Fig. 1c) separated Eastern Pacific freshwater individuals exclusively from the  
404 remaining samples, a pattern that is not expected by chance alone (permutation test  $P <$   
405 0.001, Supplementary Fig. 2, Supplementary Table 3). Rather, this reflects a shared  
406 adaptive response among Eastern Pacific freshwater populations. The exception was two  
407 freshwater individuals from the Eastern Pacific (ALA; Alaska) that did not group with the

408 other freshwater individuals from the Eastern Pacific but instead with Atlantic (marine and  
409 freshwater), Western Pacific (marine and freshwater) and the marine individuals from  
410 Eastern Pacific (Fig. 1c, Supplementary Fig. 5). LD-cluster 29 (2,728 loci, 0.109% of the  
411 dataset) covering a known inversion on Chr. XXI (Fig. 1d) grouped the Eastern Pacific  
412 freshwater individuals (except the two Alaskan individuals above) together with six Atlantic  
413 freshwater individuals and one Eastern Pacific marine individual. Because this LD-cluster  
414 maps to an inversion, the groups also represent putative inversion karyotypes. Thus, this  
415 inversion shows strong ecotype differentiation not only in the Eastern Pacific, but also in a  
416 small proportion of individuals outside of the Eastern Pacific that putatively carry the  
417 freshwater-adapted karyotype (i.e. the karyotype with the highest frequency among  
418 Eastern Pacific freshwater individuals). Notably, no cluster of similar magnitude to LD-  
419 cluster 2 – which separates freshwater individuals from one specific region from all  
420 remaining samples in the data – could be detected outside of the Eastern Pacific,  
421 demonstrating that parallel marine-freshwater differentiation in the Eastern Pacific is much  
422 more prevalent than anywhere else in the world.

423 A small proportion of the loci from LD-cluster 2 (28 SNPs) mapped to regions that showed  
424 global parallelism in Jones et al. (2012; Supplementary Table 3). In addition, <1% of all  
425 loci in LD-cluster 2 (243 SNPs) showed  $F_{ST} > 0.2$  also in the Atlantic (as is evident e.g.  
426 from Fig. 2a and Supplementary Fig. 4a). These loci appear to be non-randomly  
427 distributed in the genome (Supplementary Fig. 4a), indicating that indeed they are likely to  
428 be linked to genomic regions involved in marine-freshwater differentiation in both the  
429 Atlantic and the Eastern Pacific. Due to small sample-sizes, the  $F_{ST}$  Manhattan plots  
430 display a considerable amount of noise, particularly in the datasets from the Eastern and  
431 Western Pacific Oceans (Fig. 2b,e and Supplementary Fig. 4).

432

### 433 *Geographic structure and regional local adaptation*

434 LD-cluster 1 (10,184 loci, 0.405% of the dataset) separated all Pacific individuals (Eastern  
435 and Western) from the Atlantic individuals (Fig. 1a), thus mainly reflecting trans-oceanic  
436 geographic structure. LD-clusters 4, 8, 9, 14 and 24 (a total of 526 loci, 0.021% of the  
437 dataset, Supplementary Fig. 2) separated freshwater individuals from only one geographic  
438 region; this likely reflects geographic clustering but could also contain some loci involved in  
439 non-parallel freshwater adaptation. These loci are therefore interpreted as inconclusive  
440 with respect to their underlying evolutionary phenomena (Supplementary Table 3).  
441 Accordingly, loci from these LD-clusters showed little marine-freshwater differentiation in  
442 both the Eastern Pacific and Atlantic (Supplementary Fig. 2), and only 2 loci (from LD-  
443 cluster 14) mapped to the *m-f* global regions (Supplementary Table 3). All freshwater  
444 individuals in this dataset were important, as they can inform us about the geographic  
445 scale of parallel marine-freshwater differentiation (for the LD-clusters where marine-  
446 freshwater differentiation also involved geographic regions where marine samples were  
447 available).

448

### 449 **Proof of concept simulations**

450 In the simulated data, before *Atl* was colonized from *Pac*, all five freshwater populations in  
451 *Pac* were fixed or nearly fixed for the freshwater-adapted alleles of all locally adapted QTL  
452 (Fig. 3f). Following the colonization of *Atl* (38-440 Kya; Fig. 3f), the increased frequency of  
453 the freshwater allele in the Atlantic freshwater populations depended on both QTL density  
454 and the level of gene flow between *Pac* and *Atl* (Fig. 3f). The highest increase in  
455 freshwater-adapted alleles in *Atl* was observed when QTL density was low (3 QTL per  
456 chromosome) and trans-oceanic gene flow was high (5 migrants/generation, Fig. 3f).

457 During post-glacial colonization of new freshwater habitats from the sea (10 Kya to  
458 present), freshwater-adapted alleles (in both *Pac* and *Atl*) gradually increased in the newly  
459 formed freshwater populations (Fig. 3f), reflecting local adaptation. This increase was  
460 similarly dependent on the QTL density (both in *Pac* and *Atl*) and trans-oceanic gene flow  
461 (only affecting *Atl*, Fig. 3f). These patterns likely reflect the underlying levels of ancestral  
462 variation in the sea available for subsequent freshwater adaptation (Supplementary Fig.  
463 6a). The lowest frequencies of freshwater-adapted alleles in the sea were always  
464 observed when QTL density was the highest (in both *Pac* and *Atl*) and trans-oceanic gene  
465 flow was the lowest (only affecting the *Atl*, Supplementary Fig. 6a). Furthermore, the  
466 frequency of freshwater-adapted alleles in both the sea (ancestral variation) and in the  
467 post-glacial freshwater populations (local adaptation) depended on whether the QTL were  
468 located in low or high recombination regions; the lowest frequencies of freshwater-adapted  
469 alleles were always observed in low recombination regions (Supplementary Fig. 6b,c). The  
470 freshwater-adapted alleles in both *Pac* and *Atl* freshwater populations never reached  
471 similar frequencies during the post-glacial colonization (10 Kya; Fig. 3f) as before post-  
472 glacial colonization (>10 Kya; Fig. 3f), showing that ancestral variation in the sea was not  
473 sufficient to allow complete local adaptation (i.e. fixation of all original freshwater-adapted  
474 alleles) in our simulations. Note that with the rapid fixation of all the freshwater-adapted  
475 alleles (that started at frequency 0.5 in the freshwater populations in *Pac*) and the low  
476 mutation rate used ( $1e^{-8}$  per site and generation), the contribution of *de novo* mutations (at  
477 the QTL) to freshwater adaptation in these simulations are negligible.

478 In the simulations, present-day marine-freshwater differentiation (mean neutral  $F_{ST}$ ) was  
479 always low for the two chromosomes without QTL, and high recombination regions of  
480 chromosomes that contain QTL (Fig. 4; Supplementary Fig. 6d). In contrast,  $F_{ST}$  for low  
481 recombination regions of QTL-containing chromosomes was high for *Pac* (for all

482 parameter settings), indicating strong islands of parallel marine-freshwater differentiation.  
483 This was also true for *Atl* when QTL density was low (3 or 6 QTL per chromosome) and  
484 when trans-oceanic gene flow was high, but not when QTL density was high (9 QTL per  
485 chromosome; Fig. 4; Supplementary Fig. 6d) and trans-oceanic gene flow was low.

486 In the LD-clusters with the strongest *PF* versus non-*PF* differentiation in the simulated  
487 data, CSSs were always high (>0.75) between *EF* and the Atlantic populations (*AF* and  
488 *AM*), similar to LD-cluster 2 (Fig. 3g-i). However, in contrast to the simulated data, the CSS  
489 between *EF* and *PM* for LD-cluster 2 was also high (0.62), whereas for the simulated data  
490 this score increased with QTL density (starting from < 0.2), and more so when migration  
491 rate during colonization of *Atl* from *Pac* was high (5 migrants/generation; Fig. 3g-i). This is  
492 likely due to QTL density increasing the distance between *PF* and *PM* and migration rate  
493 decreasing the distance between *Pac* and *Atl* individuals, as CSS here is scaled by the  
494 maximum Euclidean distance between any two points in the data. Furthermore, when  
495 migration rate was low, no LD-cluster showed any significant CSS between *AF* and *AM*.  
496 However, when migration rate was high and with increasing QTL density, LD-clusters  
497 similar to LD-cluster 2 were readily observed also in *Atl*. This shows that in simulations,  
498 low migration rates and high QTL densities are required to produce patterns similar the  
499 observed data.

500

## 501 Discussion

502 Using genome-wide SNP data from a comprehensive global sampling of marine and  
503 freshwater stickleback ecotypes, we demonstrate that a much smaller proportion of the  
504 genome (0.208% of the dataset) is involved in global parallel marine-freshwater  
505 differentiation than exclusively in the Eastern Pacific (2.149% of the dataset). This shows  
506 that parallel evolution in the three-spined stickleback is much more pervasive in the

507 Eastern Pacific than anywhere else in the world. Indeed, the LD signal from marine-  
508 freshwater differentiation in the Eastern Pacific is even stronger than that from geographic  
509 structuring between the Pacific and Atlantic Oceans – LD-clusters separating freshwater  
510 individuals from the Eastern Pacific comprised five times as many loci than the LD-cluster  
511 reflecting geographic structuring between Pacific and Atlantic Oceans. With simulations,  
512 we demonstrate that this pattern could partly be explained by the stochastic loss of low  
513 frequency freshwater-adapted alleles in the sea during range expansion from the Eastern  
514 Pacific. As predicted, the discrepancy between the simulated Pacific and Atlantic  
515 populations in both  $F_{ST}$  and CSS analyses was the highest when trans-oceanic gene flow  
516 was low (stronger founder event), but this also required the QTL density of locally adapted  
517 loci to be high, as this reduced the levels of standing variation of freshwater-adapted  
518 alleles in the Atlantic. However, the loss of ancestral variation due to founder effects and  
519 the transporter hypothesis is likely not the only explanation for the large discrepancy in  
520 patterns of marine-freshwater differentiation between the Eastern Pacific and Atlantic  
521 Oceans. In the following, we discuss alternative biological processes that could potentially  
522 contribute to this discrepancy.

523

#### 524 *Geographic heterogeneity in standing genetic variation*

525 The “transporter hypothesis”(Schluter & Conte 2009) postulates that a low frequency of  
526 freshwater-adapted alleles is maintained in the sea via recurrent gene flow between  
527 ancestral marine and previously colonized freshwater populations. This standing genetic  
528 variation is what selection acts on during the subsequent colonization of freshwater  
529 habitats. This implicitly assumes that a large pool of locally adapted alleles has  
530 accumulated over a long period of time, as gene flow is expected to spread potentially  
531 beneficial mutations across demographically independent populations (Johannesson *et al.*

532 2010; Kempainen *et al.* 2011). In support of this hypothesis, it has been shown that  
533 haplotypes repeatedly used in freshwater adaptation are identical by descent (Colosimo *et*  
534 *al.* 2005; Roesti *et al.* 2014) and old – on average, six million years (My), but some are  
535 reported to be as old as 15 My (Nelson & Cresko 2018). Notably, these studies analyzed  
536 populations from the Eastern Pacific region, which represents the oldest and most  
537 ancestral marine population (Fang *et al.* 2020; Fang *et al.* 2018) where three-spined  
538 sticklebacks are thought to have persisted since the split from their close relative, the nine-  
539 spined stickleback (*Pungitius pungitius*), approximately 26 Mya (Betancur *et al.* 2015;  
540 Matschiner *et al.* 2011; Meynard *et al.* 2012; Sanciangco *et al.* 2016; Varadharajan *et al.*  
541 2019). However, populations in the Western Pacific and the Atlantic are much younger, as  
542 they were colonized from the Eastern Pacific during the late Pleistocene (36.9-346.5 kya  
543 (Fang *et al.* 2019; Fang *et al.* 2020). Furthermore, there is no evidence for trans-oceanic  
544 admixture (Fang *et al.* 2019; Fang *et al.* 2020) following the split of Pacific and Atlantic  
545 clades, and there are no extant populations of three-spined sticklebacks in arctic Russia  
546 between the Kara Sea and the Eastern Siberian Sea. Thus, the spread of freshwater-  
547 adapted alleles from the Eastern Pacific to elsewhere via migration through the Bering  
548 Strait is unlikely, and has probably not occurred in recent times. Our simulations show that  
549 following colonization of freshwater populations from the sea, the accessibility of  
550 freshwater-adapted alleles – which is a function of colonization history, QTL-density and  
551 recombination rate – largely determines the number of loci that show high marine-  
552 freshwater differentiation. Thus, consistent with previous simulations (Feder & Nosil 2010;  
553 Roesti *et al.* 2014), genomic islands of differentiation in linked neutral loci require several  
554 QTL to cluster in low recombination regions (Fig. 4 and Supplementary Fig. 6d).  
555 Furthermore, when trans-Atlantic gene flow was low and QTL density was high, we readily

556 observed LD-clusters that showed high marine-freshwater differentiation only in the  
557 Eastern Pacific, not in the Atlantic.

558 Our simulations and empirical data suggest that both stochastic loss of genetic diversity  
559 and selection against freshwater-adapted variants likely played a role in reducing the pool  
560 of standing genetic variation of freshwater-adapted alleles in the Atlantic region. During  
561 range expansions, genetic diversity is expected to decrease with distance from the source  
562 population from which the expansion started (Ramachandran *et al.* 2005). This pattern  
563 was very clear in our simulations as well as in the empirical data, both of which show a  
564 very clear and statistically significant reduction in heterozygosity in the Atlantic region as  
565 compared to the Eastern Pacific region (Fig. 3i, Supplementary Information 3). These  
566 results are consistent with some level of founder effect following colonisation of the Atlantic  
567 basin from the Eastern Pacific (Supplementary Information 3), which could account for the  
568 random loss of standing genetic variation of freshwater adapted alleles. Furthermore,  
569 since freshwater-adapted alleles are selected against in the sea and thus occur at low  
570 frequencies in marine environments, they are even less likely to spread to new geographic  
571 regions than neutral alleles (Halliburton & Halliburton 2004; Hyten *et al.* 2006). Consistent  
572 with this, the mean individual heterozygosity for LD-cluster 2 loci was 29 times higher in  
573 the Pacific compared to the Atlantic Ocean (Supplementary Fig. 7b): a very pronounced  
574 difference compared to that in the rest of the genome.

575 Alternative explanations for the observed discrepancy in patterns of marine-freshwater  
576 differentiation between the Eastern Pacific and Atlantic include *i*) stronger spatial genetic  
577 structure in marine populations outside of the Eastern Pacific causing heterogeneity in  
578 standing genetic variation available for freshwater adaptation, and *ii*) heterogeneity in  
579 selective regimes among freshwater habitats, both between Atlantic and Eastern Pacific  
580 Oceans and between different geographic areas in the Atlantic. We have further tested



581 these hypotheses but found little or inconclusive support in our data and in other studies  
582 (Supplementary Information 3).

583

#### 584 *Secondary contact in the Eastern Pacific*

585 All of the above hypotheses assume that the original source of the ancestral variation in  
586 the Eastern Pacific and elsewhere is the same. That is, ancient Eastern Pacific marine  
587 populations carried most ancestral variation of freshwater adapted alleles at low  
588 frequencies, sourcing the Atlantic region with freshwater adapted alleles which were  
589 partially lost either stochastically or due to selection. However, an alternative hypothesis is  
590 that modern Eastern Pacific freshwater variants were not present in the marine ancestors  
591 but rather in land-locked, ice-lake freshwater populations (Bierne et al .2013). As glaciers  
592 melted, those populations could have followed meltwater downstream, establishing  
593 freshwater populations with a different stock of alleles. Secondary contact with marine  
594 sticklebacks during this time might have eroded genetic differentiation across most of the  
595 genome, with the exception of those regions involved in freshwater adaptation. In this  
596 scenario, the standing genetic variation responsible for Eastern Pacific freshwater  
597 adaptation may not have entered the Eastern Pacific marine population until after the end  
598 of the last glaciation (i.e. after the closing of the Bearing Strait), with no potential for gene  
599 flow to the Atlantic.

600 There is ample evidence for large ice-lakes during the last glacial period (LGP) in North  
601 America, with little (if any) connection to the sea (Baker & Bunker 1985; Bretz 1969; Oviatt  
602 2015; Upham 1896). Thus, a large part of the genetic variation underlying marine and  
603 freshwater adaptation in the Eastern Pacific could in principle have evolved in allopatry i.e.  
604 separately among the freshwater ice-lake populations and in the sea (and any other

605 potential freshwater water bodies the sea is in contact with). Since the existence of these  
606 ice-lakes preceded the invasion of the Atlantic Ocean (Baker & Bunker 1985; Oviatt 2015)  
607 the freshwater-adapted alleles potentially residing in such ice-lakes could not have easily  
608 spread to the Atlantic. Consistent with this hypothesis is the strong pattern of long-range  
609 LD observed among Eastern Pacific marine individuals (Hohenlohe *et al.* 2012), as well as  
610 our LDna results which revealed one large cluster separating the Pacific and Atlantic  
611 Ocean individuals, and one that specifically separates all Eastern Pacific freshwater  
612 individuals from all other individuals. The secondary contact hypothesis is also consistent  
613 with close to zero  $H_E$  among Atlantic marine individuals observed for LD-cluster 2 loci  
614 (Supplementary Fig. 7a). Curiously, we found significant isolation by distance in the  
615 Atlantic but not in the Eastern Pacific where overall population structuring was  
616 nevertheless higher than in the Atlantic (Supplementary Fig. 7d). This could be consistent  
617 with the secondary contact hypothesis, if introgression was stronger in some regions of the  
618 Eastern Pacific compared to others. However, further empirical and simulation studies are  
619 needed to test the extent to which this secondary contact hypothesis provides a better  
620 explanation for the observed data than the transporter hypothesis alone.

621

### 622 *Conditions that allow global parallelism*

623 Genomic islands of parallel ecotype divergence were more likely to arise in the simulations  
624 when several QTL clustered in the same low recombination region. Surprisingly these  
625 were also the QTL where the frequency of the freshwater-adapted allele showed the  
626 lowest frequencies in the sea and thus, were least likely to spread to Atlantic during  
627 colonisation from Pacific. Since QTL in low recombination regions are less likely to be  
628 separated by recombination when freshwater-adapted individuals migrate to the sea, it is  
629 reasonable to assume that the selection pressure against these “haplotypes” in the sea is

630 stronger (Hohenlohe et al., 2012). However, this is not consistent with the empirical data  
631 showing that the genomic regions most likely to show global parallel ecotype divergence  
632 are inversions, where recombination in heterokaryotypes is particularly restricted. Our  
633 simulations assume that freshwater-adapted alleles are selected against in the sea (and  
634 the strength of this selection is equal for all QTL) while in reality, selection against some of  
635 the "freshwater haplotypes/karyotypes" in the sea may be weak or even absent, allowing  
636 them to easily spread during range expansions. Consistent with this reasoning, in PCAs  
637 based on loci from LD-clusters corresponding to inversions (LD-clusters 6, 22 and 29)  
638 several marine individuals also cluster with the freshwater individuals (Fig. 1d,e,  
639 Supplementary Fig. 2), indicating frequent occurrence of the "freshwater karyotypes" in the  
640 sea. Indeed, Terekhanova et al. (2019) found that the genomic regions most commonly  
641 involved in local adaptation in multiple independent freshwater populations were also  
642 those with the highest frequencies in the sea. In other words, the most geographically  
643 widespread genomic regions involved in freshwater adaptation (*sensu* the transporter  
644 hypothesis) are likely to experience the weakest selection against them in the sea,  
645 allowing them to remain at higher frequencies in the sea as standing genetic variation  
646 (Terekhanova *et al.* 2019).

647

648 *Are three-spined sticklebacks a representative model to study parallel evolution?*

649 Since the pattern of parallel genetic differentiation between marine and freshwater  
650 stickleback ecotypes in the Eastern Pacific is in stark contrast to what is seen across other  
651 parts of the species distribution range, it is reasonable to question the generality of the  
652 findings from the Eastern Pacific stickleback studies with respect to parallel evolution on  
653 broader geographic scales. A recent review of parallel evolution suggests that even  
654 dramatic phenotypic parallelism can be generated by a continuum of parallelism at the

655 genetic level (Bolnick et al., 2018). For instance, the coastal ecotypes of *Senecio laetus*  
656 exhibit only partial reuse of particular QTL among replicate populations (Roda et al., 2017),  
657 and genetic redundancy frequently underlies polygenic adaptation in *Drosophila* (Barghi et  
658 al., 2019). Similarly, using  $F_{ST}$  outliers to detect putative genomic targets of selection,  
659 Kautt et al. (2012, cichlid fishes), Le Moan et al. (2016, anchovy) and Westram et al.  
660 (2014, periwinkles) showed that phenotypically very similar populations often share only a  
661 small proportion of their  $F_{ST}$  outliers.

662 One exception that seems more general across taxa is the repeated involvement of  
663 chromosomal inversions in parallel evolution. Chromosomal inversions could store  
664 standing variation as a balanced polymorphism and distribute it to fuel parallel adaptation  
665 (Morales et al., 2019). For instance, the same Chr. I inversion involved in global marine-  
666 freshwater differentiation in three-spined sticklebacks (Jones et al., 2012, Terekhanova et  
667 al., 2014, 2019, Liu et al., 2018, this study) also differentiates stream and lake ecotypes in  
668 the Lake Constance basin in Central Europe (Roesti et al. 2015). Two other clear  
669 examples where most genetic differentiation between ecotypes at larger geographic scales  
670 is partitioned into inversions come from monkey flowers (*Mimulus guttatus*; Twyford and  
671 Friedman, 2015) and marine periwinkles (*Littorina saxatilis*; Faria et al., 2018; Westram et  
672 al., 2018).

673 While our study focuses exclusively on marine-freshwater ecotype pairs of three-spined  
674 sticklebacks, other ecotype pairs within freshwater habitats, such as stream vs. lake and  
675 benthic vs. limnetic, also exist. A recent study focusing on stream-lake populations found  
676 that putative selected loci showed greater parallelism in the Eastern Pacific (Vancouver  
677 Island) than the global scale (North America and Europe; Paccard et al., 2020), i.e. a  
678 similar pattern as reported by our study. Furthermore, Conte et al. (2015) studied the  
679 extent of QTL reuse in parallel phenotypic divergence of limnetic and benthic three-spined

680 sticklebacks within Paxton and Priest Lakes (British Columbia), and found that although  
681 76% of 42 phenotypic traits diverged in the same direction, only 49% of the underlying  
682 QTL evolved in parallel in both lakes. For highly parallel traits in two other pairs of benthic-  
683 limnetic sticklebacks, only 32% of the underlying QTL were reported to be shared (Conte  
684 *et al.* 2012). Thus, these studies are in stark contrast to the original conclusions of  
685 widespread genetic parallelism in three-spined sticklebacks. Notably, the two freshwater  
686 individuals from the Eastern Pacific that did not cluster with the remaining freshwater  
687 individuals from the Eastern Pacific (and were subsequently removed from the datasets  
688 used for  $F_{ST}$  genome scans) were from Alaska. These two individuals are also  
689 phylogenetically distinct from other freshwater individuals from the Eastern Pacific (Fang *et*  
690 *al.* 2018). One explanation for this could be that the highly divergent freshwater  
691 populations in the Eastern Pacific have a different colonization history than the Alaskan  
692 lakes. More specifically, the former could have been colonized from some divergent ice-  
693 lake refugia (see above), whereas the latter could have independently been colonized from  
694 the sea.

695

## 696 **Conclusions**

697 Our results demonstrate that genetic parallelism in the marine-freshwater three-spined  
698 stickleback model system is in fact not as pervasive as some earlier studies focusing on  
699 Eastern Pacific populations have led us to believe. Our analysis of geographically more  
700 comprehensive data, with similar and less assumption-burdened methods as used in  
701 earlier studies, shows that the extraordinary genetic parallelism observed in the Eastern  
702 Pacific Ocean is not detectable elsewhere in the world (e.g. Atlantic Ocean, Western  
703 Pacific Ocean). Hence, the focus on the Eastern Pacific has generated a perception bias –

704 the patterns detected there do not actually apply to the rest of the world. Furthermore, our  
705 simulations show that the spread of freshwater-adapted alleles can be hampered if  
706 colonization of the Atlantic from the Pacific was limited, particularly for QTL clustered in  
707 low recombination regions (i.e. those most likely to result in parallel islands of ecotype  
708 differentiation). Therefore, geographic differences in the incidence and pervasiveness of  
709 parallel evolution in three-spined sticklebacks likely stem from geographic heterogeneity in  
710 access to, and amount of, standing genetic variation, which in turn has been influenced by  
711 selection as well as historical population demography. Such historical demographic factors  
712 include founder events as well as the potential accumulation of genetic ecotypic  
713 differences in allopatry during the last glacial maximum, followed by a secondary contact  
714 only after the Atlantic Ocean was colonized via the sea from the Eastern Pacific. Hence,  
715 while striking genome-wide patterns of genetic parallelism exist (e.g. in Eastern Pacific  
716 sticklebacks), the conditions under which such patterns can occur may be far from  
717 common, perhaps even exceptional.

718

## 719 **Acknowledgements**

720 We are grateful to the following people who helped in obtaining the samples used in this  
721 study: Jacquelin DeFaveri, Anders Adill, Windsor Aguirre, Theo Bakker, Alison Bell, Mike  
722 Bell, Bertil Borg, Fredrik Franzén, Akira Goto, Andrew Hendry, Gabor Herczeg, Frank von  
723 Hippel, Aki Hirvonen, Jenni Hämäläinen, Markku Kaukoranta, Agnieszka Kijewska, David  
724 Kingsley, Yoshinobu Kosaka, Lotta Kvarnemo, Dmitry Lajus, Tuomas Leinonen, Arne  
725 Levsen, Scott McCairns, Antoine Millet, Jola Morozinska, Corey Munk, Hannu Mäkinen,  
726 Arne Nolte, Kjartan Østbye, Wäinö Pekkola, Jouko Pokela, Mark Ravinet, Katja Räsänen,  
727 Dolph Schluter, Mat Seymor, Takahito Shikano, Per Sjöstrand, Garrett Staines, Björn

728 Stelbrink, Ilkka Syv nper , Anti Vasem gi, Mike Webster, James Willacker, Helmut  
729 Winkler, and Linda Zaveik. Our research was supported by grants from Academy of  
730 Finland (250435, 263722, 265211 and 1307943 to JM and grant 316294 to PM), the  
731 Finnish Cultural Foundation (grant 00190489 to PK) and the Chinese Scholarship Council  
732 (grant 201606270188 to BF). We wish to thank three anonymous referees for constructive  
733 criticism, and Jacquelin De Faveri for feedback and linguistic corrections.

734

### 735 **Author contributions**

736 PK and JM conceive the concept of the study, with developments from PM and BF. BF, PK  
737 and PM performed analyses. PK, BF and PM wrote the manuscript. XF contributed to  
738 lifeover analysis. BF visualised the data. JM contributed to shaping the research and  
739 manuscript. All authors accepted the final version.

## References

- 740  
741  
742 Arendt J, Reznick D (2008) Convergence and parallelism reconsidered: what have we learned  
743 about the genetics of adaptation? *Trends in Ecology & Evolution* **23**, 26-32.
- 744 Baker VR, Bunker RC (1985) Cataclysmic Late Pleistocene Flooding from Glacial Lake Missoula -  
745 a Review. *Quaternary Science Reviews* **4**, 1-41.
- 746 Barghi N, Tobler R, Nolte V, *et al.* (2019) Genetic redundancy fuels polygenic adaptation in  
747 *Drosophila*. *PLoS Biology* **17**, e3000128.
- 748 Bell MA, Foster SA (1994) *The evolutionary biology of the threespine stickleback* Oxford University  
749 Press.
- 750 Betancur RR, Orti G, Pyron RA (2015) Fossil-based comparative analyses reveal ancient marine  
751 ancestry erased by extinction in ray-finned fishes. *Ecology Letters* **18**, 441-450.
- 752 Bolnick DI, Barrett RDH, Oke KB, Rennison DJ, Stuart YE (2018) (Non)Parallel Evolution. *Annual*  
753 *Review of Ecology Evolution and Systematics* **49**, 303-330.
- 754 Bretz JH (1969) The Lake Missoula floods and the channeled scabland. *Journal of Geology* **77**,  
755 505-543.
- 756 Catchen J, Hohenlohe PA, Bassham S, Amores A, Cresko WA (2013) Stacks: an analysis tool set  
757 for population genomics. *Molecular Ecology* **22**, 3124-3140.
- 758 Chan YF, Marks ME, Jones FC, *et al.* (2010) Adaptive evolution of pelvic reduction in sticklebacks  
759 by recurrent deletion of a Pitx1 enhancer. *Science* **327**, 302-305.
- 760 Colosimo PF, Hosemann KE, Balabhadra S, *et al.* (2005) Widespread parallel evolution in  
761 sticklebacks by repeated fixation of Ectodysplasin alleles. *Science* **307**, 1928-1933.
- 762 Conte GL, Arnegard ME, Best J, *et al.* (2015) Extent of QTL Reuse During Repeated Phenotypic  
763 Divergence of Sympatric Threespine Stickleback. *Genetics* **201**, 1189-1200.
- 764 Conte GL, Arnegard ME, Peichel CL, Schluter D (2012) The probability of genetic parallelism and  
765 convergence in natural populations. *Proceedings: Biological sciences, The Royal Society*  
766 **279**, 5039-5047.
- 767 DeFaveri J, Shikano T, Shimada Y, Goto A, Merila J (2011) Global analysis of genes involved in  
768 freshwater adaptation in threespine sticklebacks (*Gasterosteus aculeatus*). *Evolution* **65**,  
769 1800-1807.
- 770 Fang B, Kempainen P, Momigliano P, Merilä J (2019) Oceans apart: Heterogeneous patterns of  
771 parallel evolution in sticklebacks. *BioRxiv*, 826412.



- 772 Fang B, Merila J, Matschiner M, Momigliano P (2020) Estimating uncertainty in divergence times  
773 among three-spined stickleback clades using the multispecies coalescent. *Molecular*  
774 *Phylogenetics and Evolution* **142**, 106646.
- 775 Fang B, Merila J, Ribeiro F, Alexandre CM, Momigliano P (2018) Worldwide phylogeny of three-  
776 spined sticklebacks. *Molecular Phylogenetics and Evolution* **127**, 613-625.
- 777 Faria R, Chaube P, Morales HE, *et al.* (2018) Multiple chromosomal rearrangements in a hybrid  
778 zone between *Littorina saxatilis* ecotypes. *Molecular Ecology* **28**, 1375-1393.
- 779 Feder JL, Nosil P (2010) The efficacy of divergence hitchhiking in generating genomic islands  
780 during ecological speciation. *Evolution* **64**, 1729-1747.
- 781 Ferchaud AL, Hansen MM (2016) The impact of selection, gene flow and demographic history on  
782 heterogeneous genomic divergence: three-spine sticklebacks in divergent environments.  
783 *Molecular Ecology* **25**, 238-259.
- 784 Fox EA, Wright AE, Fumagalli M, Vieira FG (2019) ngsLD: evaluating linkage disequilibrium using  
785 genotype likelihoods. *Bioinformatics* **35**, 3855-3856.
- 786 Gibson G (2005) The synthesis and evolution of a supermodel. *Science* **307**, 1890-1891.
- 787 Halliburton R, Halliburton R (2004) *Introduction to population genetics* Pearson/Prentice Hall Upper  
788 Saddle River, NJ.
- 789 Hedrick PW (2007) Sex: differences in mutation, recombination, selection, gene flow, and genetic  
790 drift. *Evolution* **61**, 2750-2771.
- 791 Hendry AP, Peichel CL, Matthews B, Boughman JW, Nosil P (2013) Stickleback research: the now  
792 and the next. *Evolutionary Ecology Research* **15**, 111-141.
- 793 Hohenlohe PA, Bassham S, Currey M, Cresko WA (2012) Extensive linkage disequilibrium and  
794 parallel adaptive divergence across threespine stickleback genomes. *Philos Trans R Soc*  
795 *Lond B Biol Sci* **367**, 395-408.
- 796 Hohenlohe PA, Bassham S, Etter PD, *et al.* (2010) Population genomics of parallel adaptation in  
797 threespine stickleback using sequenced RAD tags. *PLoS Genetics* **6**, e1000862.
- 798 Hohenlohe PA, Magalhaes IS (2019) The population genomics of parallel adaptation: lessons from  
799 threespine stickleback. In: *Population Genomics*. Springer. Cham.
- 800 Hu A, Meehl GA, Otto-Bliesner BL, *et al.* (2010) Influence of Bering Strait flow and North Atlantic  
801 circulation on glacial sea-level changes. *Nature Geoscience* **3**, 118-121.
- 802 Hubbard T, Andrews D, Cáccamo M, *et al.* (2005) Ensembl 2005. *Nucleic Acids Research* **33**,  
803 D447-D453.

- 804 Hyten DL, Song Q, Zhu Y, *et al.* (2006) Impacts of genetic bottlenecks on soybean genome  
805 diversity. *Proceedings of the National Academy of Sciences of the United States of America*  
806 **103**, 16666-16671.
- 807 Johannesson K, Panova M, Kempainen P, *et al.* (2010) Repeated evolution of reproductive  
808 isolation in a marine snail: unveiling mechanisms of speciation. *Philos Trans R Soc Lond B*  
809 *Biol Sci* **365**, 1735-1747.
- 810 Jones FC, Grabherr MG, Chan YF, *et al.* (2012) The genomic basis of adaptive evolution in  
811 threespine sticklebacks. *Nature* **484**, 55-61.
- 812 Kautt AF, Elmer KR, Meyer A (2012) Genomic signatures of divergent selection and speciation  
813 patterns in a 'natural experiment', the young parallel radiations of N icaraguan crater lake  
814 cichlid fishes. *Molecular Ecology* **21**, 4770-4786.
- 815 Kempainen P, Knight CG, Sarma DK, *et al.* (2015) Linkage disequilibrium network analysis  
816 (LDna) gives a global view of chromosomal inversions, local adaptation and geographic  
817 structure. *Molecular Ecology Resources* **15**, 1031-1045.
- 818 Kempainen P, Lindskog T, Butlin R, Johannesson K (2011) Intron sequences of arginine kinase in  
819 an intertidal snail suggest an ecotype-specific selective sweep and a gene duplication.  
820 *Heredity* **106**, 808-816.
- 821 Kitano J, Ross JA, Mori S, *et al.* (2009) A role for a neo-sex chromosome in stickleback speciation.  
822 *Nature* **461**, 1079.
- 823 Korneliussen TS, Albrechtsen A, Nielsen R (2014) ANGSD: Analysis of Next Generation  
824 Sequencing Data. *BMC Bioinformatics* **15**, 356.
- 825 Le Moan A, Gagnaire PA, Bonhomme F (2016) Parallel genetic divergence among coastal–marine  
826 ecotype pairs of European anchovy explained by differential introgression after secondary  
827 contact. *Molecular Ecology* **25**, 3187-3202.
- 828 Lescak EA, Bassham SL, Catchen J, *et al.* (2015) Evolution of stickleback in 50 years on  
829 earthquake-uplifted islands. *Proceedings of the National Academy of Sciences of the*  
830 *United States of America* **112**, E7204-7212.
- 831 Li H (2011) A statistical framework for SNP calling, mutation discovery, association mapping and  
832 population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987-  
833 2993.
- 834 Li H, Durbin R (2009) Fast and accurate short read alignment with Burrows–Wheeler transform.  
835 *Bioinformatics* **25**, 1754-1760.
- 836 Li Z, Kempainen P, Rastas P, Merila J (2018) Linkage disequilibrium clustering-based approach  
837 for association mapping with tightly linked genomewide data. *Molecular Ecology Resources*  
838 **18**, 809-824.

- 839 Liu S, Ferchaud AL, Gronkjaer P, Nygaard R, Hansen MM (2018) Genomic parallelism and lack  
840 thereof in contrasting systems of three-spined sticklebacks. *Molecular Ecology* **27**, 4725-  
841 4743.
- 842 Matschiner M, Hanel R, Salzburger W (2011) On the origin and trigger of the notothenioid adaptive  
843 radiation. *PLoS ONE* **6**, e18911.
- 844 Matthey-Doret R, Whitlock MC (2019) Background selection and FST: consequences for detecting  
845 local adaptation. *Molecular Ecology* **28**, 3902-3914.
- 846 McKinnon JS, Rundle HD (2002) Speciation in nature: the threespine stickleback model systems.  
847 *Trends in Ecology & Evolution* **17**, 480-488.
- 848 Meiri M, Lister AM, Collins MJ, *et al.* (2014) Faunal record identifies Bering isthmus conditions as  
849 constraint to end-Pleistocene migration to the New World. *Proceedings: Biological*  
850 *sciences, The Royal Society* **281**, 20132167.
- 851 Meynard CN, Mouillot D, Mouquet N, Douzery EJ (2012) A phylogenetic perspective on the  
852 evolution of Mediterranean teleost fishes. *PLoS ONE* **7**, e36443.
- 853 Natri HM, Shikano T, Merilä J (2013) Progressive recombination suppression and differentiation in  
854 recently evolved neo-sex chromosomes. *Molecular Biology and Evolution* **30**, 1131-1144.
- 855 Nelson TC, Cresko WA (2018) Ancient genomic variation underlies repeated ecological adaptation  
856 in young stickleback populations. *Evolution letters* **2**, 9-21.
- 857 Neuenschwander S, Hospital F, Guillaume F, Goudet J (2008) quantiNemo: an individual-based  
858 program to simulate quantitative traits with explicit genetic architecture in a dynamic  
859 metapopulation. *Bioinformatics* **24**, 1552-1553.
- 860 Orti G, Bell MA, Reimchen TE, Meyer A (1994) Global survey of mitochondrial DNA sequences in  
861 the threespine stickleback: evidence for recent migrations. *Evolution* **48**, 608-622.
- 862 Oviatt CG (2015) Chronology of Lake Bonneville, 30,000 to 10,000 yr BP. *Quaternary Science*  
863 *Reviews* **110**, 166-171.
- 864 Paccard A, Hanson D, Stuart YE, *et al.* (2020) Repeatability of Adaptive Radiation Depends on  
865 Spatial Scale: Regional Versus Global Replicates of Stickleback in Lake Versus Stream  
866 Habitats. *Journal of Heredity* **111**, 43-56.
- 867 Pujolar JM, Ferchaud AL, Bekkevold D, Hansen MM (2017) Non-parallel divergence across  
868 freshwater and marine three-spined stickleback *Gasterosteus aculeatus* populations.  
869 *Journal of Fish Biology* **91**, 175-194.
- 870 Ramachandran S, Deshpande O, Roseman CC, *et al.* (2005) Support from the relationship of  
871 genetic and geographic distance in human populations for a serial founder effect originating  
872 in Africa. *Proceedings of the National Academy of Sciences of the United States of America*  
873 **102**, 15942-15947.

- 874 Roda F, Walter GM, Nipper R, Ortiz-Barrientos D (2017) Genomic clustering of adaptive loci during  
875 parallel evolution of an Australian wildflower. *Molecular Ecology* **26**, 3687-3699.
- 876 Roesti M, Gavrillets S, Hendry AP, Salzburger W, Berner D (2014) The genomic signature of  
877 parallel adaptation from shared genetic variation. *Molecular Ecology* **23**, 3944-3956.
- 878 Roesti M, Kueng B, Moser D, Berner D (2015) The genomics of ecological vicariance in threespine  
879 stickleback fish. *Nature Communications* **6**, 8767.
- 880 Roesti M, Moser D, Berner D (2013) Recombination in the threespine stickleback genome--  
881 patterns and consequences. *Molecular Ecology* **22**, 3014-3027.
- 882 Sanciangco MD, Carpenter KE, Betancur RR (2016) Phylogenetic placement of enigmatic  
883 percomorph families (Teleostei: Percomorphaceae). *Molecular Phylogenetics and Evolution*  
884 **94**, 565-576.
- 885 Schaffner SF (2004) The X chromosome in population genetics. *Nature Reviews Genetics* **5**, 43-  
886 51.
- 887 Schluter D, Conte GL (2009) Genetics and ecological speciation. *Proceedings of the National*  
888 *Academy of Sciences of the United States of America* **106 Suppl 1**, 9955-9962.
- 889 Stankowski S, Chase MA, Fuiten AM, *et al.* (2019) Widespread selection and gene flow shape the  
890 genomic landscape during a radiation of monkeyflowers. *PLoS Biology* **17**, e3000391.
- 891 Stern DL (2013) The genetic causes of convergent evolution. *Nature Reviews Genetics* **14**, 751-  
892 764.
- 893 Terekhanova NV, Barmintseva AE, Kondrashov AS, Bazykin GA, Mugue NS (2019) Architecture of  
894 Parallel Adaptation in Ten Lacustrine Threespine Stickleback Populations from the White  
895 Sea Area. *Genome Biology and Evolution* **11**, 2605-2618.
- 896 Terekhanova NV, Logacheva MD, Penin AA, *et al.* (2014) Fast evolution from precast bricks:  
897 genomics of young freshwater populations of threespine stickleback *Gasterosteus*  
898 *aculeatus*. *PLoS Genetics* **10**, e1004696.
- 899 Twyford AD, Friedman J (2015) Adaptive divergence in the monkey flower *Mimulus guttatus* is  
900 maintained by a chromosomal inversion. *Evolution* **69**, 1476-1486.
- 901 Upham W (1896) *The glacial lake agassiz* US Government Printing Office.
- 902 Varadharajan S, Rastas P, Loytynoja A, *et al.* (2019) A high-quality assembly of the nine-spined  
903 stickleback (*Pungitius pungitius*) genome. *Genome Biology and Evolution*.
- 904 Westram A, Galindo J, Alm Rosenblad M, *et al.* (2014) Do the same genes underlie parallel  
905 phenotypic divergence in different *L. ittorina saxatilis* populations? *Molecular Ecology* **23**,  
906 4603-4616.

907 Westram AM, Rafajlović M, Chaube P, *et al.* (2018) Clines on the seashore: The genomic  
908 architecture underlying rapid divergence in the face of gene flow. *Evolution letters* **2**, 297-  
909 309.

910 Zheng X, Levine D, Shen J, *et al.* (2012) A high-performance computing toolset for relatedness  
911 and principal component analysis of SNP data. *Bioinformatics* **28**, 3326-3328.

912 Östlund-Nilsson S, Mayer I, Huntingford FA (2006) *Biology of the three-spined stickleback* CRC  
913 Press.

914

915

916 **FIGURES**

---

917 **Figure 1 | Linkage Disequilibrium network analysis (LDna).** (a-h) Eight main clusters of  
918 loci identified by LDna (LD-clusters). In each panel (LD-cluster), the top and middle plots  
919 present the marine-freshwater differentiation ( $F_{ST}$ ) between Atlantic and Eastern Pacific  
920 samples, respectively. The bottom plot shows the principal component analysis (PCA)  
921 based on the LD-cluster loci. The seven different colours represent the geographic origin  
922 of populations. Solid and open circles refer to freshwater and marine ecotypes,  
923 respectively. All identified LD-clusters (29 in total) and corresponding information are  
924 presented in Supplementary Fig. 2 and Supplementary Table 2. (j) Map of the sampled  
925 populations; colours match those in the PCA results. A Mercator projection of the sampling  
926 map is shown in Supplementary Fig. 1.

927

928 **Figure 2 | Genetic parallelism identified by the unsupervised and supervised**  
929 **methods.** (a) Comparison of marine-freshwater differentiation ( $F_{ST}$ ) in the Atlantic (x-axis)  
930 and Eastern Pacific (y-axis) datasets for the three LD-clusters (LD-clusters 2, 21 and 29)  
931 associated with strong marine-freshwater parallelism in the Eastern Pacific. (b) Genome-  
932 wide  $F_{ST}$  of the Eastern Pacific samples for loci of the LD-clusters coloured as in (a). (c)  
933 The same as (a) but for the twelve LD-clusters (5, 6, 10, 11, 12, 13, 16, 18, 20, 22, 25 and  
934 27) that are involved in global marine-freshwater genetic parallelism. (d) and (e) Genome-  
935 wide  $F_{ST}$  of the Atlantic and Eastern Pacific samples, respectively, with colours  
936 corresponding to LD-cluster loci in (c). The position of the Ectodysplasin (EDA) locus is  
937 indicated in (d) and (e).

938

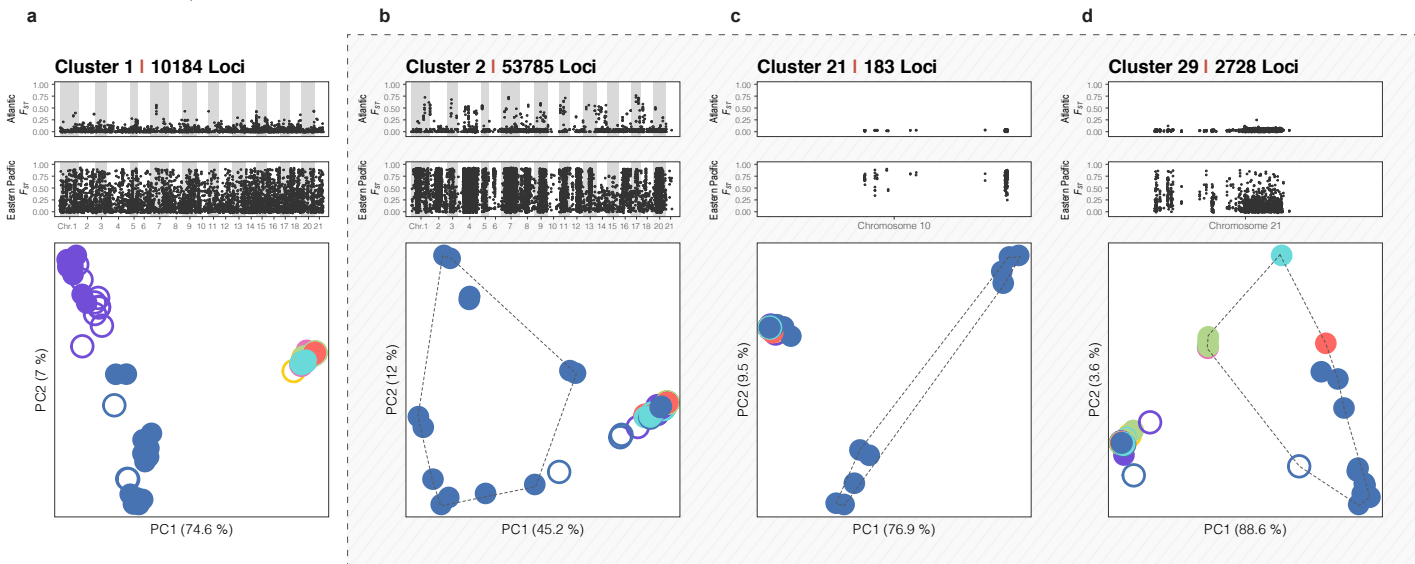
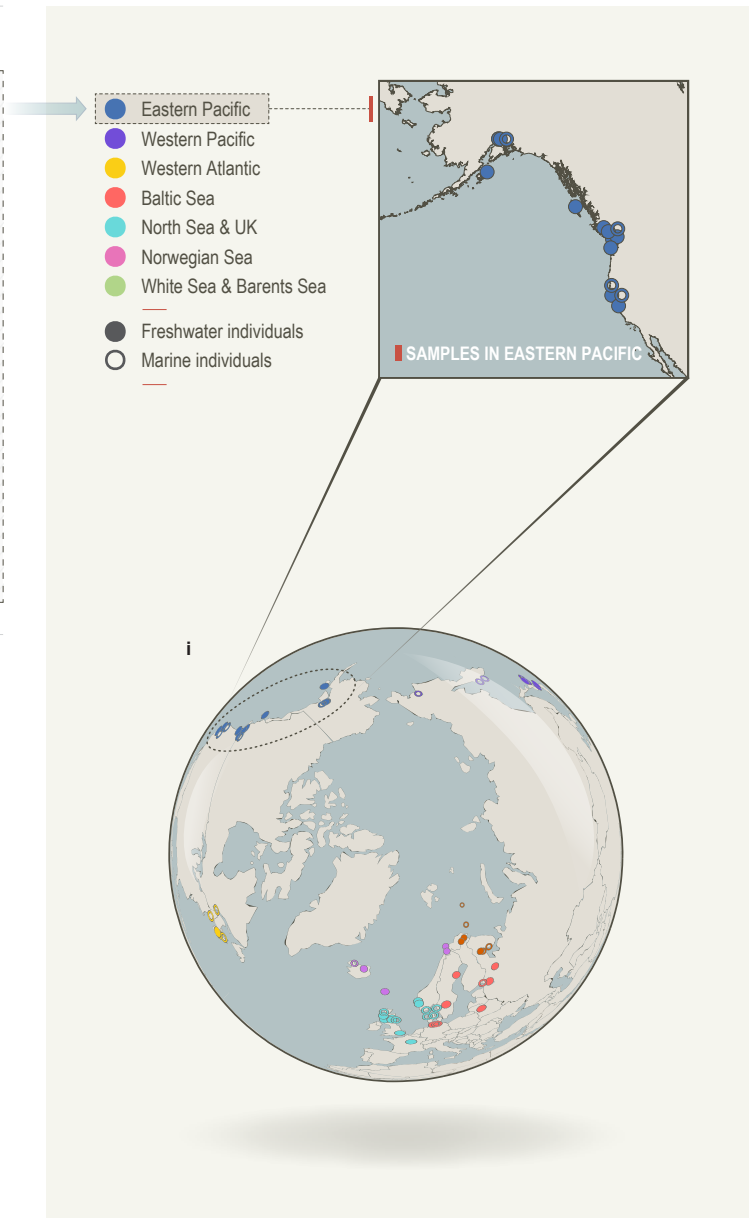
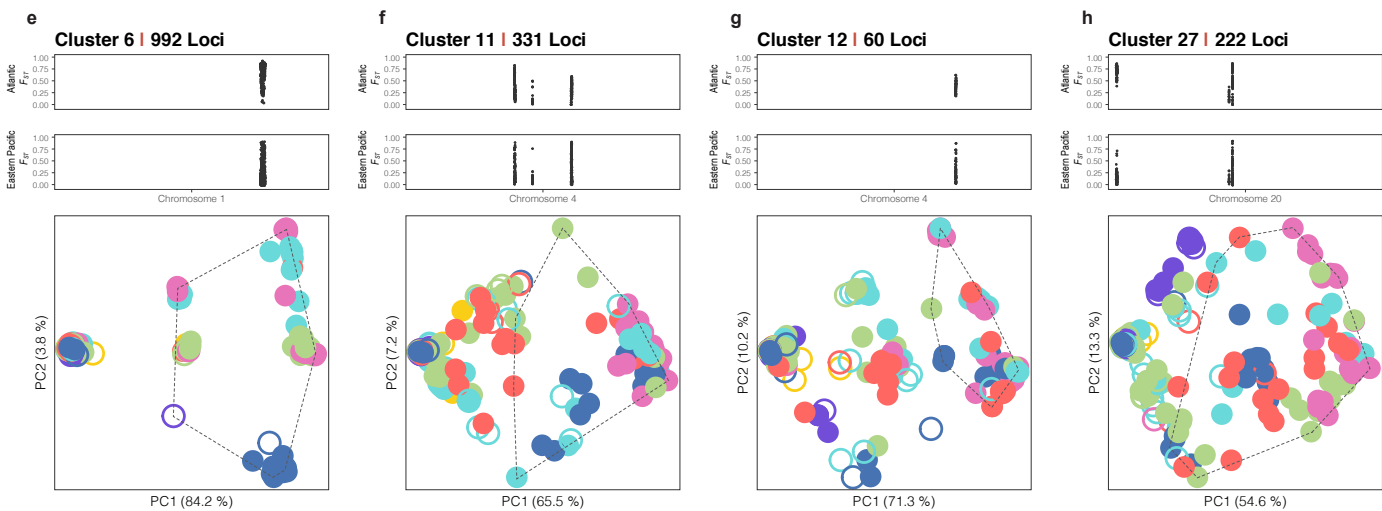
939 **Figure 3 | Ecological genetics in simulated data.** (a-e) A schematic of the demographic  
940 scenario used for the simulations that is consistent with the “transporter hypothesis” of  
941 Schluter & Conte (2009). (a) Initial local adaption of the freshwater populations in the  
942 Pacific. (b) The colonization of stickleback populations from the Pacific to the Atlantic. (c)  
943 Geographic isolation between the two oceans. (d) Extinction of lakes during the last glacial  
944 period (LGP) with the survival of refuge populations. (e) The post-glacial colonization of  
945 the new freshwater populations. (f) Frequency of selected (freshwater-adapted) alleles in  
946 the newly established freshwater populations through generations at high and low levels of  
947 trans-oceanic gene flow and different QTL-densities. (g) PCA of the empirical data (LD-  
948 cluster 2; left) and the simulated data (right), with ecotypes and geographical regions as  
949 shown in the figure legend. (h) Cluster separation score (CSS) of the empirical and  
950 simulated data in the Pacific and Atlantic oceans, respectively. (i) Observed heterozygosity  
951 in different geographical regions in the empirical and simulated data (Supplementary  
952 Information 3).

953

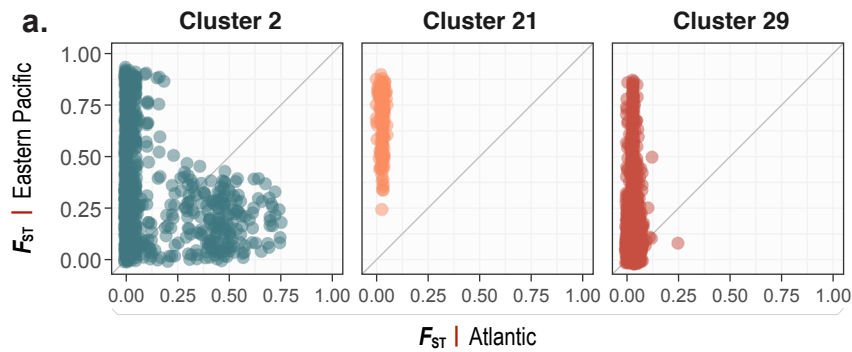
954 **Figure 4 | Genomic differentiation in simulated data.** (a, b) Genome-wide marine-  
955 freshwater differentiation ( $F_{ST}$ ) from simulated data (data from the last generation  
956 representing present day sampling). For each parameter combination, loci from all 20  
957 replicates were pooled. Red dots indicate QTL, and blue dots indicate loci from LD-  
958 clusters that were the most similar to LD-cluster 2 (empirical data) showing the strongest  
959 marine-freshwater differentiation in the Eastern Pacific (grey represent non-LD cluster  
960 loci). (c)  $F_{ST}$  distribution of QTL in the simulations (all replicates pooled), indicating the  
961 proportion of loci that were either fixed ( $F_{ST} \sim 1$ ), lost ( $F_{ST} \sim 0$ ), or were fixed to different  
962 degrees in only 1, 2 or 3 of the four freshwater populations ( $0.1 \lesssim F_{ST} \lesssim 0.9$ ) since post

963 glacial colonisation. A small amount of noise (along the x-axis) has been added to the QTL  
964 positions to improve their visibility.

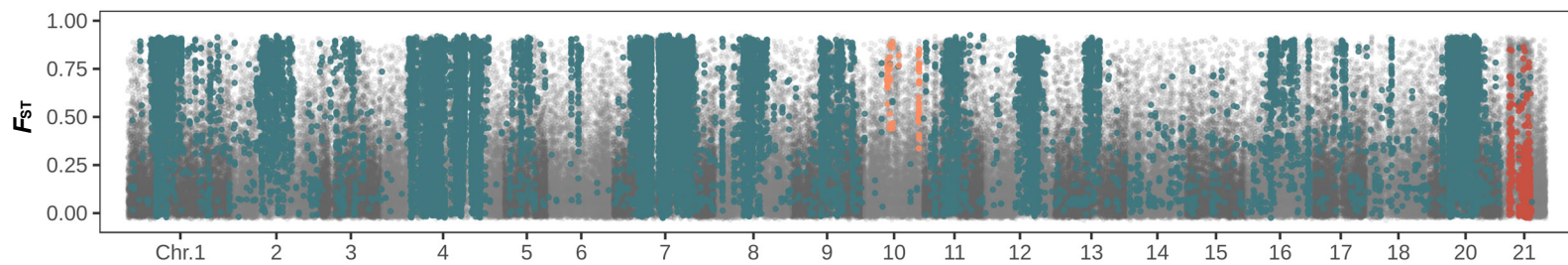


**FIGURE 1****GEOGRAPHIC STRUCTURE****GENETIC PARALLELISM • EASTERN PACIFIC (EXCLUSIVELY AND DOMINANTLY)****GENETIC PARALLELISM • TRANS-OCEANIC (REPRESENTATIVES)**

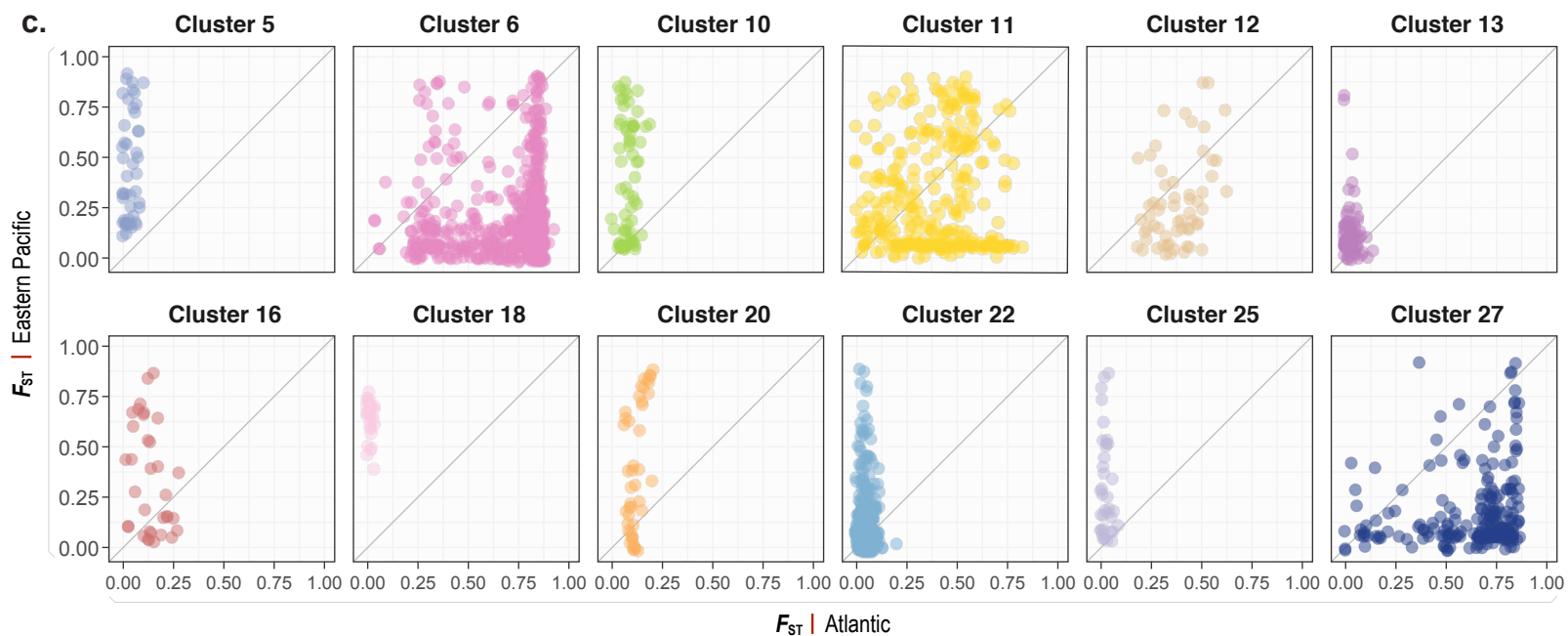
GENETIC PARALLELISM • EASTERN PACIFIC (EXCLUSIVE AND DOMINANT)



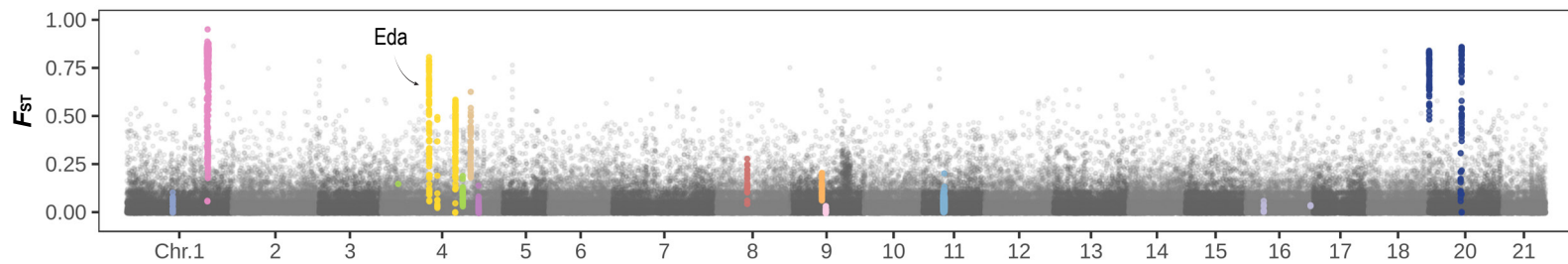
**b.** MARINE-FRESHWATER DIFFERENTIATION IN THE EASTERN PACIFIC



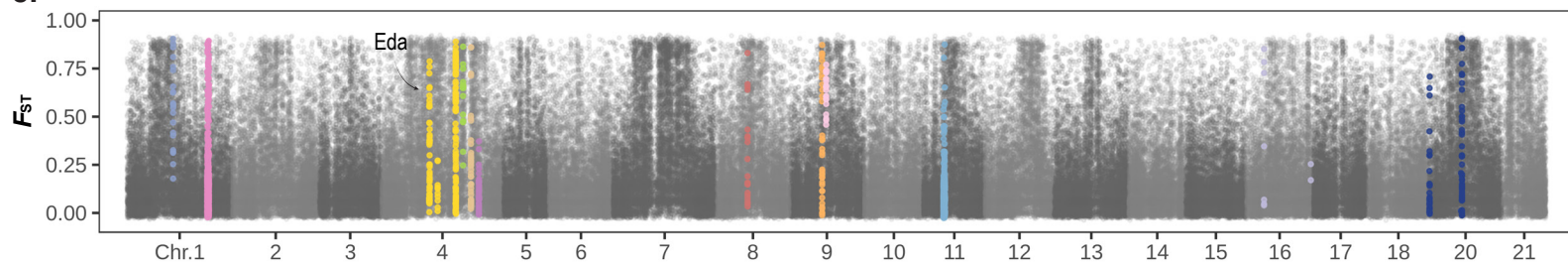
GENETIC PARALLELISM • TRANS-OCEANIC



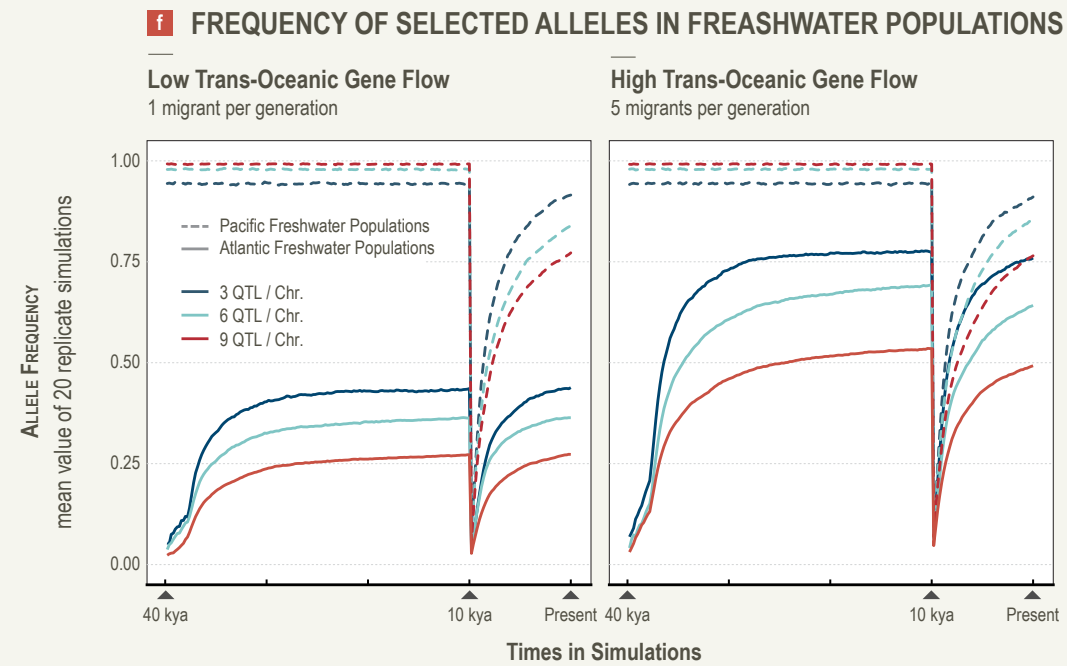
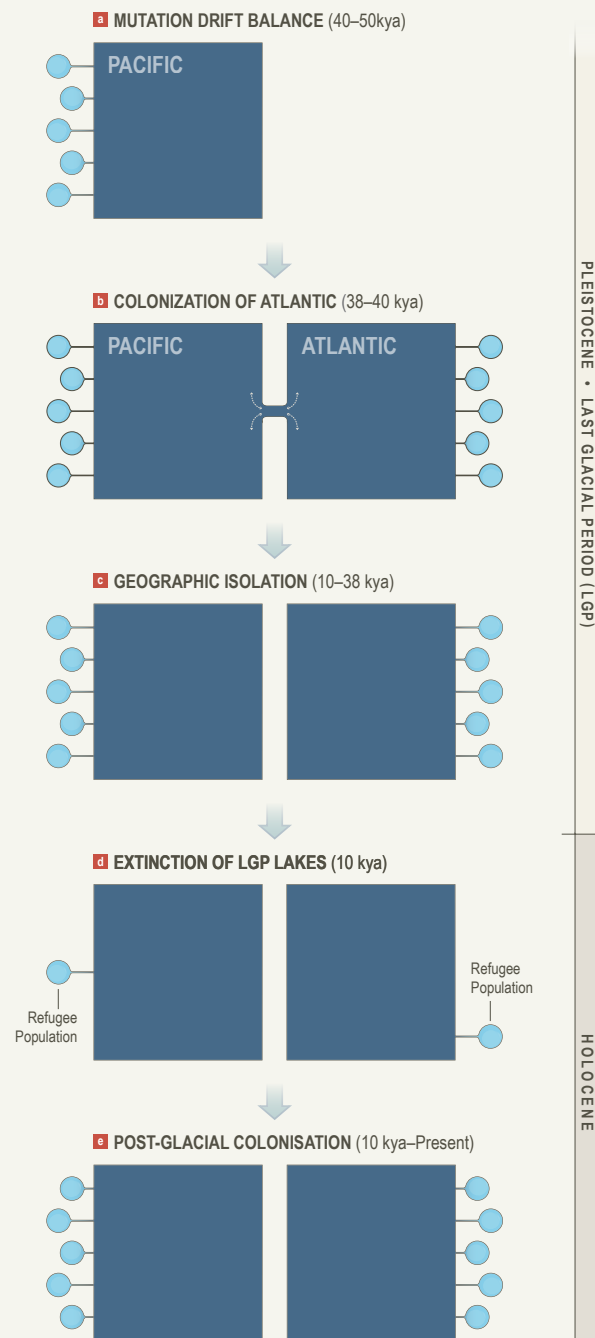
**d.** MARINE-FRESHWATER DIFFERENTIATION IN THE ATLANTIC



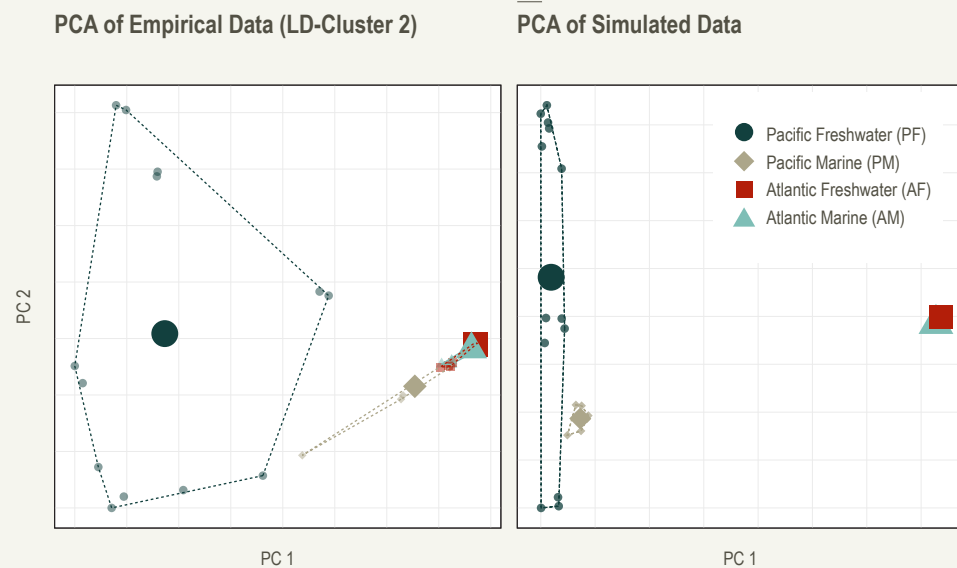
**e.** MARINE-FRESHWATER DIFFERENTIATION IN THE EASTERN PACIFIC



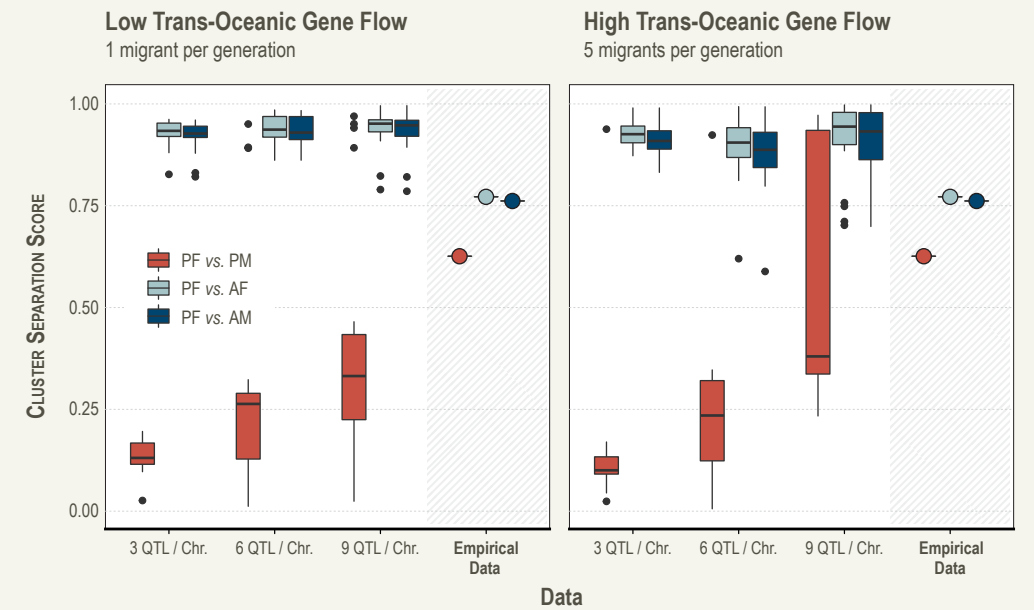
**FIGURE 3**



**g** PCA OF EMPIRICAL AND SIMULATED DATA



**h** CLUSTER SEPARATION SCORE IN THE PACIFIC OCEAN



**i** HETEROZYGOSITY

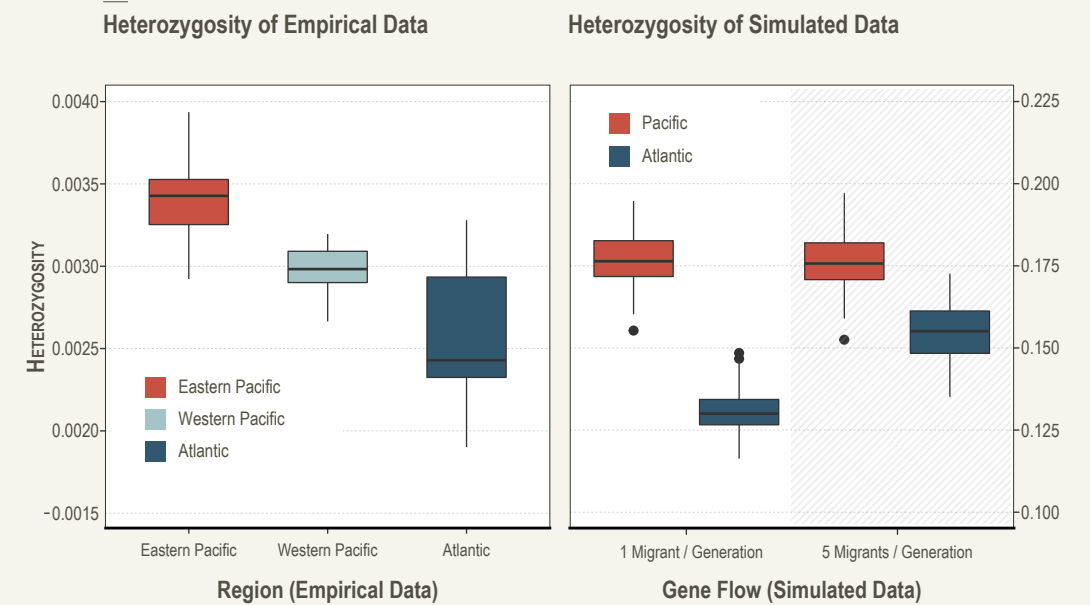


FIGURE 4  
a.

