# Origin Matters: Using a Local Reference Genome Improves Measures in Population Genomics

Doko-Miles J. Thorburn[1,2+], Kostas Sagonas[1,3], Mahesh Binzer-Panchal[4], Frederic J.J. Chain[5], Philine G.D. Feulner[6,7], Erich Bornberg-Bauer[8], Thorsten BH Reusch[9], Irene E. Samonte-Padilla[10], Manfred Milinski[10], Tobias L. Lenz[11,12], Christophe Eizaguirre[1]

[1]School of Biological and Chemical Sciences, Queen Mary University of London, London, United Kingdom.

[2]Department of Life Sciences, Imperial College London, London, United Kingdom.

[3]Department of Zoology, School of Biology, Aristotle University of Thessaloniki, Panepistimioupoli, 54124 Thessaloniki, Greece.

[4]Department of Medical Biochemistry and Microbiology, National Bioinformatics Infrastructure Sweden (NBIS), Science for Life Laboratory, Uppsala University, Sweden.

[5]Department of Biological Sciences, University of Massachusetts Lowell, USA.

[6]Department of Fish Ecology and Evolution, Centre of Ecology, Evolution and Biogeochemistry, EAWAG Swiss Federal Institute of Aquatic Science and Technology, Kastanienbaum, Switzerland.

[7]Division of Aquatic Ecology and Evolution, Institute of Ecology and Evolution, University of Bern, Bern, Switzerland.

[8]Evolutionary Bioinformatics, Institute for Evolution and Biodiversity, Westfälische Wilhelms University, Münster, Germany.

[9]Marine Evolutionary Ecology, GEOMAR Helmholtz Centre for Ocean Research, Kiel, Germany.

[10]Department of Evolutionary Ecology, Max Planck Institute for Evolutionary Biology, Plön, Germany.

[11]Research Group for Evolutionary Immunogenomics, Max Planck Institute for Evolutionary Biology, Plön, Germany.

[12]Research Unit for Evolutionary Immunogenomics, Department of Biology, University of Hamburg, Hamburg, Germany.

+Corresponding author: d.m.j.thorburn@qmul.ac.uk

# Abstract

Genome-level sequencing enables us to ask fundamental questions about the genetic basis of adaptation, population structure, and epigenetic mechanisms, but usually requires a suitable reference genome for mapping population-level re-sequencing data. In some model systems, multiple reference genomes are available, giving researchers the challenging task of determining which reference genome best suits their data. Here we compare the use of two different reference genomes for the three-spined stickleback (*Gasterosteus aculeatus*), one novel genome derived from a European gynogenetic individual and the published reference genome of a North American individual. Specifically, we investigate the impact of using a local reference versus one generated from a distinct lineage on several common population genomics analyses. Through mapping genome resequencing data of 60 sticklebacks from across Europe and North America, we demonstrate that genetic distance among samples and the reference impacts downstream analyses. Using a local reference genome increased mapping efficiency and genotyping accuracy, effectively retaining more and better data. Despite comparable distributions of the metrics generated across the genome using SNP data (i.e., $\pi$, Tajima's $D$, and $F_{ST}$), window-based statistics using different references resulted in different outlier genes and enriched gene functions. A marker-based analysis of DNA methylation distributions had a comparably high overlap in outlier genes and functions, yet with distinct differences depending on the reference genome. Overall, our results highlight how using a local reference genome decreases reference bias to increase confidence in downstream analyses of the data. Such results have significant implications in all reference-genome-based population genomic analyses.

**Keywords:** Reference genomes, population genomics, gynogenetic, genome assembly, read mapping, *Gasterosteus aculeatus*, stickleback, reference mapping bias

## Introduction

Genome-level sequencing has revolutionized many biological fields including evolution, ecology, microbiology, and population genomics (M. R. Jones & Good, 2016; Kao et al., 2014; Stapley et al., 2010). Historically, scientists have relied on one or a small number of high-quality reference genomes to address their specific questions. Progressively, the availability of high-quality reference genomes from large-scale projects (e.g., Earth BioGenome and Vertebrate Genome Projects; Lewin et al., 2018; Scientists, 2009), the decreasing costs of sequencing, and the availability of curated variant databases (e.g., dbSNP and dbVar; Lappalainen et al., 2013; Sherry et al., 2001) have improved the breadth and depth of genomic research. However, limitations due to the availability and quality of these resources and the challenge of integrating and interpreting multiple sources of genetic variation still exist (Formenti et al., 2022). These problems are only exacerbated in non-model systems where a high-quality reference genome may not exist. Recently, a graph-based genome alignment methodology has been developed to index and incorporate variant databases, providing a practical and highly efficient method that better captures variation in a genome (Kim et al., 2019). However, databases of high-quality variants are not currently available for most study systems. Moreover, in highly variable regions graph-based approaches can incur significant computational overhead and have an increased change of false-positive alignments (Grytten et al., 2020; Pritt et al., 2018). Hence, until graph-based methods become common place, assembling a haploid reference is still a viable alternative to making and utilising a resource intensive genome graph.

Reference genomes are integral parts of the analytical frameworks of genomic research. Specifically, variants are identified and referred to by their loci in relation to their position mapped on a reference genome. Such an approach allows for direct comparisons among multiple individuals and organisms, therein enabling comprehensive research in various fields, including phylogenomics (Fang et al., 2018; Wang et al., 2013), comparative evolution (Palmer

83  & Kronforst, 2020; Parker et al., 2013), adaptive radiations (Ronco et al., 2020), and the

84  genomics of domestication (Frantz et al., 2015).

85  Arguably one of the most important factors to consider when multiple reference genomes or

86  assembly versions are available is their difference in quality. Whilst we are moving towards

87  complete and error-free assemblies (Rhie et al., 2021), the continuing advances of

88  methodologies can create significant differences among assemblies. For example,

89  bioinformatic tools have been developed to resolve false gene duplications that stem from

90  heterozygosity in homologous haplotypes (Guan et al., 2020; Roach et al., 2018). In some

91  cases, such as in humans, reference genomes are continually updated alongside major

92  advances, where choosing the most updated version will offer the most accurate analysis of

93  human sequencing data (Guo et al., 2017; International Human Genome Sequencing

94  Consortium, 2001, 2004; Pushkarev et al., 2009). However, when multiple similar quality and

95  genetically diverse reference genomes are available from multiple populations or strains (e.g.,

96  Berner et al., 2019; Gan et al., 2011; Hirsch et al., 2016; Springer et al., 2018), genetic distance

97  among samples and the reference may be an important factor to consider.

98  The detection of genomic polymorphisms is affected by the evolutionary time between the

99  individuals being sequenced and the reference genome (Bohling, 2020; Prasad et al., 2022;

100  B. N. Reid et al., 2021). This has implications on both the detection and the genotyping of

101  single nucleotide polymorphisms (SNPs) and structural variants (SVs). For example, variant

102  calling through pipelines such as GATK and freebayes use Bayesian inference to call

103  genotypes (i.e., the likelihood of a genotype, given the data; Auwera et al., 2014; Garrison &

104  Marth, 2012). When high differentiation between reference genome and the sampled

105  individuals exists, a high proportion of segregating sites will emerge as fixed differences

106  between the samples and the reference genome (i.e., homozygote non-reference in diploids),

107  and therefore uncertain genotypes can have a higher likelihood of being called as homozygote

108  non-reference, even if this is not the case. Analogously, most methods used to detect SVs

109  (read-depth, split paired-read, breakpoints, and assembly) compare mapped reads to the

110   reference genome (Pirooznia et al., 2015). Several studies on humans have demonstrated

111   that ethnicity-specific reference genomes are beneficial (Ameur et al., 2017; Dewey et al.,

112   2011; Fakhro et al., 2016; Lacaze et al., 2019). Specifically, the targeted reference genomes

113   improved reliability of genetic and structural variation calls (Ameur et al., 2017; Fakhro et al.,

114   2016). Moreover, recent studies have demonstrated that increasing phylogenetic distance

115   between target species and reference genome decreases mapping efficiency and has strong

116   effects on evolutionary inferences made from the data (Bohling, 2020; Prasad et al., 2022).

117   The three-spined stickleback (*Gasterosteus aculeatus*) is a supermodel in evolutionary biology

118   (K. Reid et al., 2021), and research on this small teleost fish has pioneered discoveries related

119   to the genomics of adaptation (Feulner et al., 2015; Haenel et al., 2019; F. C. Jones et al.,

120   2012; Roesti et al., 2015), adaptive divergence (Feulner et al., 2015; Huang et al., 2016; Roesti

121   et al., 2013) and development (Shapiro et al., 2004; Spitz et al., 2001). Genomic investigations

122   using *G. aculeatus* have mostly relied on a single high-quality chromosome-level reference

123   genome from an isolated population in Alaska (F. C. Jones et al., 2012), which has been

124   updated with multiple improvements (Glazer et al., 2015; Nath et al., 2021; Peichel et al., 2017;

125   Roesti et al., 2013). Recently, several additional *de novo G. aculeatus* contig-level genome

126   assemblies have been made available, including European and marine derived assemblies

127   (Berner et al., 2019). These additional genome assemblies may be appropriate given the wide

128   geographic distribution of the species; there are Atlantic and Pacific clades of *G. aculeatus*

129   that diverged an estimated 44.6 Kya with extensive phenotypic and genetic diversity across

130   its range (Fang et al., 2018; McKinnon & Rundle, 2002).

131   In this study, we report the generation and *de novo* annotation of a European *G. aculeatus*

132   genome assembly derived from a gynogen individual (Samonte-Padilla et al., 2011). The

133   gynogenesis process in *G. aculeatus* produced a near complete homozygous diploid fish

134   (Samonte-Padilla et al., 2011), which helps alleviate some of the genome assembly difficulties

135   associated with heterozygosity. Then, using European and North American derived reference

136   genomes, we investigated the effect of reference genome origin. We used high-quality

137     genome-wide resequencing data from 60 *G. aculeatus* individuals from 5 recently diverged

138     lake-river pairs distributed across the European (Atlantic clade; Fang et al., 2018) and North

139     American (Western North America; Pacific clade) *G. aculeatus* ranges. The effect of using a

140     local reference genome was investigated through comparing mapping efficiency (i.e., the

141     percentage of reads that are successfully mapped to the reference genome), genotype calling,

142     genome scans for $F_{ST}$, Tajima's *D,* and π, and identification of SVs. The resequenced genome

143     data have been used to investigate the distribution of islands of differentiation across the

144     genome (Feulner et al., 2015) and the role of copy-number variation in adaptation (Chain et

145     al., 2014), offering baselines to evaluate specific metrics against a new local genome. We also

146     compared how the different genomes affect DNA-methylation calling, focusing on an

147     additional 50 European fish for which reduced-representation bisulfite sequencing (RRBS)

148     was available (Sagonas et al., 2020). We hypothesise that the mapping and genotyping results

149     will generally fall into 3 categories: ($h_1$) a local reference genome has a significant and

150     putatively benficial effect, ($h_2$) one reference genome is better than the other, and ($h_3$) a local

151     reference genome has an effect, but the effects are more substantial with local populations

152     (i.e., a mixture of both $h_1$ and $h_2$). Our hypotheses can be observed in the outcome of statistical

153     tests: $h_1$ presents as a significant interaction between reference origin and population origin,

154     $h_2$ as a significant effect of reference origin, and $h_3$ presents similarly to $h_1$ with a significant

155     interaction, but the post-hoc analysis indicates a clear bias towards one of the reference

156     genomes.

## Material & Methods

### *Reference Genome Assembly and Annotation*

159     The induction of the diploid gynogen individual followed the protocol established for three-

160     spined sticklebacks (Samonte-Padilla et al., 2011). In brief, fertilised stickleback eggs were

161     mixed with UV-irradiated sperm (2 minutes of exposure) and then exposed to a heat-shock

162     treatment of 34°C for 4 minutes, 5 minutes post-fertilisation. The treatment caused the genetic

163  inactivation of the sperm, resulting in homozygote maternal offspring that lack paternal alleles

164  (Samonte-Padilla et al., 2011). In order to increase the likelihood of embryo development, two

165  siblings from the same family were used for this process. After 5 months post hatching, a fish

166  was sacrificed. DNA was extracted using the Qiagen high molecular weight extraction kit

167  following the manufacturer's protocol. Sequencing was then conducted at the Beijing

168  Genomics Institute (BGI), taking place on a PacBio platform. In total, 2,805,993 reads were

169  generated with a coverage of 44.1X. A total of 214,285 PacBio reads were discarded before

170  further analysis given their short length. In addition, Illumina paired-end sequencing libraries

171  in a HiSeq2500 platform were constructed with insert sizes of about 170 base pairs (bp), 500

172  bp and 800 bp.

173  Genome assembly was performed using Canu v1.6 assembler (Koren et al., 2017), followed

174  by an internal polishing step using Quiver. To hybrid polish the PacBio assembly, a total of

175  316,797,342 high-quality Illumina reads were mapped to the contigs using BWA-MEM v0.7.15

176  (Li, 2013) and the alignment was then used for further polishing using Pilon v1.22 (Walker et

177  al., 2014). Illumina raw reads were trimmed, and low quality and adaptor sequences were

178  removed using Cutadapt v1.13 (Martin, 2011). To evaluate the PacBio *de novo* contig-level

179  assembly and search for potential misassemblies, we used FRCbam v1.3.0 (Vezzi et al.,

180  2012), whereas its completeness in terms of core orthologous genes was assessed using

181  BUSCO v3.0.1 (Simão et al., 2015) and the 'actinopterygii' data set. The 32.1kb contig

182  tig00001189_pilon mapped to the mitochondrial genome of the North American *G. aculeatus*

183  (Peichel et al., 2017) and was trimmed and labelled the mitochondrial genome. The sequence

184  was trimmed to only include the 15.8kb that aligned to the North American mitochondrial

185  genome given the size of the mitochondrial genomes are moderately conserved even across

186  long phylogenetic distances (Gissi et al., 2008). The excess 16.3kb of contig

187  tig00001189_pilon were labelled and kept in the assembly, making up two additional contigs.

188  To scaffold the European *G. aculeatus* contig-level gynogen assembly into pseudo-

189  chromosomes we used Chromosemble from the Satsuma2 package (Grabherr et al., 2010).

190    Here, we ordered and oriented the contigs based on synteny with the North American *G.*

191    *aculeatus* chromosome-level genome (Peichel et al., 2017), excluding the unmapped

192    scaffolds. We tested for the effect of contig size on alignment rates, and only retained contigs

193    with an alignment rate of 70% or above for further assembly. Gynogen contigs not scaffolded

194    onto a pseudo-chromosome were concatenated in size order into an unmapped scaffold for

195    population genomic analyses, separating each contig by 1,000 N's. Synteny between the

196    European and North American *G. aculeatus* assemblies was calculated using the

197    SatsumaSynteny2 function (Grabherr et al., 2010), and plotted using the Circos v0.69

198    visualisation tool (Krzywinski et al., 2009). It should be noted that the reference genomes have

199    been created using distinctly different methodologies. As such, significant reference origin by

200    population origin interactions that do not clearly exhibit a reciprocal effect in the test statistic

201    likely include a genome quality effect.

202    Repetitive sequences were identified *de novo* in the European pseudo-chromosome assembly

203    using Repeat Modeler v.2.0.1 while Repeat Masker v.4.1.1 (Smit et al., 2015) was used to

204    mask the genome using the three-spined stickleback and zebrafish libraries in two separately

205    rounds. The results of each round were then analyzed together, complex repeats were

206    separated, to produce the final repeat annotation. Genome annotation was performed on the

207    European repeat-masked pseudo-chromosome genome assembly using MAKER2 v.2.31.9

208    (Holt & Yandell, 2011). Subsequent genome annotation was performed following a two-round

209    approach. For the first round, the repeat annotation data (release 95) as well as *G. aculeatus*

210    transcriptome and protein sequences from ENSEMBL and UniProt/SwissProt databases were

211    used as evidence sets for the prediction of gene models, while *est2genome* and

212    *protein2genome* setting were set as 1. For the second round, SNAP (Korf, 2004), with ADE of

213    0.25 and length of 50, and AUGUSTUS v.3.2.3 with default values (Stanke et al., 2008) were

214    trained on the gene model predicted from the first round. Functional annotation was performed

215    using BlastP against UniProt proteins with an E-value threshold of 1e-5, and InterProScan

216    v.5.4-47 (P. Jones et al., 2014) was used for domain annotation. The resulting gene models

217 were filtered to retain those with AED value of 0.5 or less, having PFAM annotations and

218 significant hits to known proteins against UniProt DB (E-value 1e-5).

219 We identified orthologous and paralogous gene families among the North American and

220 European *G. aculeatus* reference genomes using OrthoFinder v2.4.1 with the default

221 parameters (Emms & Kelly, 2019). Protein sequences were extracted using the getfasta

222 function in the BEDtools toolset v2.26.0, using the *-split* parameter to only include exonic

223 regions (Quinlan & Hall, 2010). Where applicable, downstream analyses were restricted to

224 only include the 7,529 1-to-1 orthologues identified between the two assemblies to remove

225 any biases stemming from differences in gene number or functional annotation.

## *SNP Data Collection and Processing*

227 Whole genome resequencing data from a total of 60 *G. aculeatus* individuals were used from

228 5 recently diverged freshwater lake (_L) and river (_R) population pairs (table S1; details on

229 sampling, library preparation, sequencing, and original data processing up to adapter trimming

230 can be found in; Feulner et al., 2015). The five population pairs were sampled from two sites

231 in Germany (G1 and G2; European; Atlantic clade), and one site in Norway (No; European;

232 Atlantic clade), the United States (US; North American; Pacific clade), and Canada (Ca; North

233 American; Pacific clade).

234 For the genome scan using SNP data, raw data was processed following previous procedures

235 (Feulner et al., 2013). Adapter-cleaned reads were trimmed using Trimmomatic v0.36 (Bolger

236 et al., 2014) in paired-end mode, trimming read tails with a PHRED quality score below 20 and

237 trimming to a maximum of 50bp. Using BWA-MEM v0.7.17 (Li, 2013), reads were

238 independently mapped to both the anchored European and US reference genomes (Peichel

239 et al., 2017). Mapping efficiency was calculated using Bamtools v2.4.1 (Barnett et al., 2011).

240 All downstream processing of mapped reads and methodology for variant calling were

241 identical for both reference genomes. Mapped reads were processed using Picard toolkit

242 v2.18.7 (https://broadinstitute.github.io/picard/), applying FixMateInformation and CleanSam.

243 All reads belonging to the same individual from different lanes were combined using

244 MergeSamFile, and then duplicate reads were flagged using MarkDuplicates. Variant calls

245 were performed using GATK v4.0.6.1 (McKenna et al., 2010), calling variants in all genomes

246 simultaneously, split by chromosome. The final set of SNPs was produced using hard filtering

247 following the best practise workflow (see supplementary methods for filtering thresholds;

248 (Depristo et al., 2011). Genome mapping and variant calls were conducted on the QMUL

249 Apocrita High Performance Computing Cluster (King et al., 2017).

250 The gynogenesis process employed here aims to purge heterozygous variation, enabling us

251 to reconstruct complex genomic regions. To assess the proportion of the genome that remains

252 heterozygotic we used the same GATK variant calling and filtering pipeline detailed above with

253 the paired-end Illumina libraries used to polish the gynogen assembly.

## *Phylogenetic Analyses*

255 In order to assess the phylogenetic relationship between the reference genomes and the 60

256 resequenced genomes, we compared maximum-likelihood phylogenies based on each variant

257 call with its respective reference genome. Maximum-likelihood phylogenies for both SNP calls

258 were inferred using RAxML v8.2.11 (Stamatakis, 2014). We randomly sampled 1% of all

259 segregating sites by setting the "—select-random-fraction" parameter to 0.01 in the GATK

260 function SelectVariants. The resulting VCFs were converted to PHYLIP format (Felsenstein,

261 1989) using the vcf2phylip.py script (doi:10.5281/zenodo.1257058). Trees were constructed

262 using the GTRGAMMA model with 1000 bootstraps. Phylogenetic trees were plotted using the

263 R package *ggtree* package v2.2.4 (Yu et al., 2017).

## *Estimations of Genotype Bias*

265 Calling of genotypes (i.e., heterozygote versus homozygote) may be affected by reference

266 genome origin. Measures of genome-wide zygosity across all segregating sites were

267 performed by parsing the genotype field in the VCF file using a custom R script. Genotypes

268 were grouped into 5 distinct categories: homozygous reference, homozygous non-reference,

269    heterozygous reference/non-reference, heterozygous non-reference/non-reference, and

270    missing. This process was repeated to estimate genotypes for each individual mapped to both

271    reference genomes.

## Genome Scan Using SNP Data

273    Genome scans were performed in R v4.0.2 (R Development Core Team, 2019), and

274    population genetics indices were calculated using the R package *PopGenome* v2.7.5 (Pfeifer

275    et al., 2014). We measured Tajima's $D$ and $\pi$ in each population, and $F_{ST}$ in each parapatric

276    population pair. All metrics were calculated in non-overlapping 20kb windows across all 20

277    autosomal chromosomes and the unmapped scaffold. To obtain outlier genomic windows, we

278    extracted the top 1% of the empirical distributions for each metric and population (or population

279    pair for $F_{ST}$). This conservative criterion (e.g., Feulner et al., 2015; Lai et al., 2019; Stern &

280    Lee, 2020) was chosen to increase our confidence in defining outliers. Finally, we defined

281    outlier genes as any gene overlapping one or more outlier genomic windows from the SNP-

282    based genome scans using the *foverlaps* function in the R package *data.table* v1.9.6 (Dowle

283    et al., 2015).

## Detection of Structural Variants

285    To investigate the impact of reference genome origin on the detection of structural variants

286    (SVs), we used two independent SV callers, DELLY2 v0.8.3 (Rausch et al., 2012) and LUMPY

287    v0.3.0 (Layer et al., 2014). Both LUMPY and DELLY were run using the default parameters.

288    The LUMPY call was genotyped using SVTyper v0.7.1 (Chiang et al., 2015), and all SVs

289    marked with the "LowQual" DELLY flag were removed. To ensure SVs found by both programs

290    were the same, only SVs with an overlap of at least 50% were accepted and merged into a

291    cohort level VCF file using SURVIVOR v1.0.3 (Jeffares et al., 2017). For downstream

292    analyses, we included autosomal duplications, deletions, and inversions with a length of

293    between 1kb and 1Mb supported by at least 6 split or discordant reads. All samples that were

294    homozygote non-reference or heterozygote across all samples were analysed separately as

295    these variants likely arise from the reference genome assemblies. Genes with coordinates

296    entirely nested within SVs were defined as structurally variable genes (SVG) using the

297    *foverlaps* function in the R package *data.table* v1.9.6.

## *Methylation Data Processing and Genome Scan*

299    For the DNA methylation analysis, we used 50 methylomes of laboratory full-sib families of

300    European *G. aculeatus* obtained from the study of Sagonas et al. (2020), were we investigated

301    whether parasite infection alters genome-wide patterns and levels of DNA methylations. For

302    each fish, a single-end library of 100 bp reads with an average of 11.5 million reads was

303    produced. The raw data were quality checked using FASTQC v0.11.5 (Andrews, 2010).

304    Cutadapt v1.9.1 (Martin, 2011) was used to trim and filter low-quality bases (-q 20), remove

305    trimmed reads shorter than 10 bases and remove adapters using multiple adapter sequences

306    (ATCGGAAGAGCACAC, AGATCGGAAGAGCACAC and NNAGATCGGAAGAGCACAC)

307    with a minimum overlap of 1 bp between adapter and read. Trimmed reads were

308    independently mapped against the European and US reference genomes (Peichel et al.,

309    2017) to extract methylated cytosines, using Bismark v0.22.1 (Krueger & Andrews, 2011) with

310    the Bowtie2 v2.3.2 aligner, allowing up to 2 mismatches. Similar to the SNP data, Bamtools

311    v2.4.1 (Barnett et al., 2011) was used to calculate the mapping efficiency.

312    Cytosine methylation ratios in CpG sites was estimated for each fish and differentially

313    methylated sites (DMS) were calculated between the two treatment groups (no parasite

314    exposed or exposed to the nematode parasite *Camallanus lacustris*) respectively (for more

315    information, see Sagonas et al., 2020) using the R package *MethylKit* v1.14.2. CpG

316    methylation ratios were estimated by calculating the number of reads mapping to a given

317    position carrying a cytosine divided by those reads carrying either C or T. CpG sites with

318    coverage below 10X and sites that have more than 99.9th percentile of coverage were

319    discarded in each sample from downstream analyses. We selected DMS between treatments

320    using the following criteria: change in fractional methylation larger than 15%; q-values lower

321    than 0.01 (SLIM method), and presence of methylated Cs in at least 50% of the samples within

322  a treatment group. Differentially methylated genes were identified using the *genomation* R

323  package v.1.1.0 (Akalin et al., 2015); a gene was considered differentially methylated if at least

324  one DMSs was located no further than 1.5-kb upstream and 500 bases downstream of it.

## *GO Enrichment Analysis*

326  To examine how the origin of the reference genome impacts functional enrichment of outlier

327  genes, structurally variable genes, and differentially methylated genes, we performed gene

328  ontology (GO) enrichment analyses. GO enrichment analyses were performed using the

329  g:GOSt function in the g:Profiler v0.1.9 package (Raudvere et al., 2019). Outlier genes were

330  grouped by all combinations of reference genome origin versus population of origin, resulting

331  in 4 distinct combinations: Europe-Europe, Europe-North America, North America-Europe,

332  and North America-North America. The lists of outlier genes or SVGs identified using the

333  European assembly were assigned the orthologous North American gene identifiers for GO

334  enrichment. *P*-values were corrected for multiple testing using FDR.

## *Statistical Analyses*

336  Linear mixed effect models in the *lme4* package (Bates et al., 2015) were used to analyse how

337  the origin of the reference genomes affected mapping efficiency, detection of genome-wise

338  zygosity, and population-genetic indices (i.e., Tajima's $D$, $\pi$, and $F_{ST}$). For mapping efficiency,

339  the proportion of mapped reads that were singletons or duplications were independently used

340  as response variables. An interaction between reference genome origin (i.e., North America

341  or Europe) and population origin (i.e., North America or Europe) were assigned as fixed

342  effects, and population ID was set as a random factor. The same approach was used for all

343  analyses, independently replacing the response variable with the zygosity categories or

344  population-genetic indices, retaining the same fixed and random effects. *P*-values were

345  inferred using Satterwaite's degree of freedom method in the *lmerTest* package (Kuznetsova

346  et al., 2017). Tukey HSD post-hoc tests were performed using *emmeans* (Lenth et al., 2020).

# Results

## *Reference Genome Assembly*

The contig-level European gynogen assembly is 458.4Mb, making it 1.0% smaller than the North American reference genome (463.0Mb; Peichel et al., 2017). We achieved an N50 of 0.746Mb, comprised of 1,906 contigs (table 1). Guided by synteny, we conservatively placed 419.3Mb (91.5%) into 21 chromosomes (fig. 1), forming a pseudochromosome level assembly allowing for a more accurate comparison of the impact of the origin of reference genomes on population genomic metrics. A combination of gene evidence from the *G. aculeatus* North American reference genome (available cDNA and protein sequences) and *ab initio* gene predictions resulted in 22,739 genes annotated and a BUSCO completeness score of 95.2%. A total of 18,255 (90.2%) genes in the European gynogen assembly were orthologous to genes in the North American assembly. Additionally, we confirmed the European and North American derived reference genomes are nested within the phylogeny of populations sampled from the same geographic regions (supplementary fig. S1).

Next, we called SNPs in the paired-end libraries used to polish the European genome assembly to assess the remaining heterozygosity in the gynogen genome assembly, which is expected to be homozygous after the gynogenesis procedure. In total we identified 299,931 SNPs (average genome-wide read depth of 60.6) in the genome: 3,118 SNPs were homozygote non-reference, and 296,813 were heterozygote. Hence, after two generations of gynogenesis only 0.07% of the genome remains polymorphic, compared to 0.53% ± 0.01% (mean ± SE) across all samples mapped to both reference assemblies (supplementary table S1). These SNPs likely stem from paralogous sequences that were not identified by the genome assembler or from errors in the SNP calling process.

## *Mapping Efficiency*

To assess the impact of reference genome origin on downstream analyses, we started by mapping whole genome resequencing reads from 60 individuals to the reference genomes. In

373    total, we identified 10,672,162 SNPs using the North American reference genome, and

374    10,757,204 SNPs using the European reference genome (supplementary table S1). Overall,

375    we achieved an average genome-wide depth of coverage of 20.40 (range; 11.09 - 35.46) using

376    the North American reference and 20.99 (11.30 - 36.14; supplementary table S1) using the

377    European reference. Reads mapped more efficiently (i.e., a higher proportion of reads

378    mapped) to the European reference genome regardless of the population of origin with total

379    reads mapping 2.31% ± 0.16% (European; estimate ± SE) and 1.13% ± 0.20% (North

380    American) more efficiently, albeit with a greater increase in efficiency for local samples (LMER;

381    Reference Origin:Population Origin, $F_{108}=21.38$, P<0.001; fig. 2A and supplementary table

382    S2). This same result was also observed when only considering properly matched paired-end

383    reads (i.e., paired reads where both reads map in the correct orientation; 0x2 flag), where

384    reads mapped more efficiently to the European reference, but there was a greater increase in

385    efficiency for local populations (European, 3.14% ± 0.22%; North American, 1.36% ± 0.27%;

386    Reference Origin:Population Origin, $F_{108}=25.94$, P<0.001; fig. 2A). There were also 0.83% ±

387    0.08% (European) and 0.24% ± 0.09% (North American) fewer singletons (i.e., where only

388    one of the paired-end reads mapped) when mapping European or North American populations

389    to the European reference (LMER; Reference Origin:Population Origin, $F_{108}=25.10$, P<0.001;

390    fig. 2B). Reference genome or population origin had no effect on mapping of duplicate reads

391    (LMER; Reference Origin:Population Origin, $F_{108}=0.01$, P=0.942; Reference Origin, $F_{108}=0.12$,

392    P=0.729; Population Origin, $F_8=3.76$, P=0.088; fig. 2B). Overall, using the European derived

393    genome assembly noticeably improved mapping efficiency irrespective of the sample origins,

394    although the effects were noticeably more efficient for European populations. Such results are

395    consistent with hypothesis $h_3$, which describes that a local reference genome is beneficial, but

396    the effects are generally more substantial when using the European reference genome.

397    *Estimations of Genotype Performance*

398    When using a local reference genome, there were 6.24% ± 0.38% (European; estimate ± SE)

399    and 0.47% ± 0.47% (North American) fewer missing genotypes (LMER; Reference

400     Origin:Population Origin, $F_{428}=123.97$, P<0.001; fig. 3 and supplementary table S2). The same

401     was true for the detection of heterozygote genotypes of both classes (i.e., heterozygote

402     reference/non-reference and non-reference/non-reference), where a local reference genome

403     decreases the number of calls by 0.13% ± 0.09% (European) and 0.31% ± 0.11% (North

404     American) for heterozygote reference/non-reference, and 0.10% ± 0.003% (European) and

405     0.05% ± 0.004% (North American) for heterozygote non-reference/non-reference genotypes

406     (LMER; all models report the Reference Origin:Population Origin interaction; reference/non-

407     reference, $F_{428}=9.86$, P=0.002; non-reference/non-reference, $F_{428}=872.85$, P<0.001; fig. 3).

408     Conversely, we identified approximately 4 times higher proportions of homozygote reference

409     genotypes when using a local reference genome for European populations (12.19% ± 0.37%)

410     in comparison to the North American populations (3.36% ± 0.35%; LMER; Reference

411     Origin:Population Origin , $F_{428}=691.30$, P<0.001; fig. 3). Finally, approximately a two-fold

412     decrease in the proportion of homozygote non-reference variants were identified when using

413     a local reference genome for European populations (-5.72% ± 0.12%) over the North American

414     populations (-2.53% ± 0.14%; LMER; Reference Origin:Population Origin , $F_{428}=2012.29$,

415     P<0.001; fig. 3). Except for duplications which were not significantly different among calls,

416     using a local reference genome offers a benefit in the detection of every class of genotype.

417     Overall, these results are consistent with hypothesis $h_1$ which posited that a local reference

418     genome would offer a benefit and would manifest as significant interactions.

419     *Genome Scan Using SNP Data*

420     The genome-wide distributions of metrics commonly used in population genomics were

421     affected by the origin of the genome used (fig. 4 & S2-3). For a genome scan investigating

422     patterns of differentiation, $F_{ST}$ was $0.019 ± 8.4×10^{-4}$ (mean ± SE) higher for European

423     populations and $0.005 ± 1.0×10^{-5}$ higher for North American populations when using a local

424     reference genome (LMER; Reference Origin:Population Origin , $F_{220053}=324.60$, P<0.001;

425     supplementary table S2). A similar pattern was observed using $T_D$, whereby $T_D$ was 0.026 ±

426     $2.95×10^{-3}$ (European) and $5.90×10^{-3} ± 3.61×10^{-3}$ (North American) higher when using a local

427    reference genome (LMER; Reference Origin:Population Origin , $F_{406709}$=63.97, P<0.001). In

428    addition, using a local reference genome led to significantly lower values of nucleotide diversity

429    being detected (LMER; Reference Origin:Population Origin , $F_{407698}$=123.20, P<0.001), with

430    low estimates for both European ($-3.28\times10^{-5} \pm 5.39\times10^{-6}$) and North American ($-9.99\times10^{-5} \pm$

431    $6.60\times10^{-6}$) populations. These results are consistent with hypothesis $h_1$, where a local

432    reference genome has a significant impact on analyses.

433    Through sampling each metric individually and taking the top ~1% of the $F_{ST}$, π, or $T_D$

434    distributions we generated multiple lists of outlier genes (table 2; supplementary table S3).

435    Notably, the distributions of outlier genes across the genome were not significantly different

436    when using either reference genome (two-sample Kolmogorov–Smirnov test; $F_{ST}$, D=0.140,

437    P=0.988; π, D=0.121, P=0.998; $T_D$, D=0.157, P=0.962; fig. 4 & S2-3). Overall, higher number

438    of outlier genes, including the subset of outlier 1-to-1 orthologous genes, were identified when

439    using a local reference genome (table 2). The difference in genes among calls with different

440    reference genomes putatively translated into no overlapping significantly enriched GO terms

441    for the $F_{ST}$, π, and $T_D$ analyses if any enrichment was detected (table 2). Complete GO

442    enrichment tables are reported in supplementary table S4.

## *Genomic Structural Variants*

444    We next addressed the question whether using a local reference genome affected the

445    detection of structural variants. Firstly, we generated multiple distributions of SVs organised

446    by the combination of reference genome and population origin (table 3). Overall, no fixed SVs

447    were observed in all samples. The distribution of deletions significantly differed among SV

448    calls (two-sample Kolmogorov–Smirnov test, D=0.278, P=0.008; fig. S4), whereby fewer SVs

449    were detected when using the North American reference genome. The distribution of

450    duplications and inversions did not significantly differ among SV calls with either reference

451    genome (two-sample Kolmogorov–Smirnov test; duplications, D=0.109, P=0.890; inversions,

452    D=0.140, P=0.754; fig. S4). Additionally, we identified less deletions when using a local

453    reference genome, more deletions were identified overall when using the European reference

454     genome (LMER; Reference Origin:Population Origin , $F_{108}$=61.813, P<0.001; table 3,

455     supplementary table S5-S6). Analogously, we identified significantly fewer inversions when

456     using a local reference genome, but more inversions were identified overall when using the

457     North American reference genome (LMER; Reference Origin:Population Origin , $F_{108}$=34.52,

458     P<0.001; table 3). Finally, a similar number of duplications were identified in European

459     populations irrespective of reference origin, whereas fewer duplications were identified in

460     North American populations when using a local reference genome (LMER, Reference

461     Origin:Population Origin, $F_{108}$=3.99, P=0.048; supplementary table S5-S6). Overall, the results

462     of the SV analysis are in line with hypothesis $h_3$ where the effects of a local reference genome

463     are present, but unequal effects indicate one reference is better than the other.

464     Next, we investigated the role of a local reference genome on the detection of genes entirely

465     nested within SVs, which we defined as structurally variable genes (SVG; table 3). This

466     analysis was limited to 1-to-1 orthologs to ensure there was no bias arising from copy number

467     variation among reference genome annotations. Here, we identified significantly fewer SVG-

468     deletions for European populations using a local reference genome, however there was no

469     significant difference in SVG-deletions in the North American populations when using either

470     reference (LMER, Reference Origin:Population Origin , $F_{108}$=25.94, P<0.001; supplementary

471     table S5-S6). On the other hand, we identified significantly more SVG-duplications when using

472     the North American reference, regardless of the origin of the population (LMER, Reference

473     Origin:Population Origin , $F_{108}$=9.40, P=0.003; supplementary table S5-S6). Finally,

474     significantly fewer inversions were detected when using a local reference genome (LMER;

475     Reference Origin:Population Origin , $F_{108}$=62.55, P<0.001; supplementary table S5-S6).

476     To identify whether a local reference genome correlated with the detection of functional

477     enrichment in SVs, we investigated the SVGs GO enrichment. Similar to the overall SVG

478     distribution analysis, this analysis was restricted to 1-to-1 orthologs. We identified no overlap

479     in functional enrichment in the majority of the comparisons (table 3). The one exception was

480     SVG-deletions in the North American populations, where 3 out of 13 significantly enriched

481 terms (signaling receptor regulator activity, signaling receptor activator activity, and receptor

482 ligand activity) were identified in both calls. Overall, the detection of SVs and SVGs were

483 affected in different ways by the origin of the reference genome.

## DNA Methylation Analysis

485 Finally, we conducted a DNA methylation analysis after mapping bisulfite sequencing reads

486 to the two reference genomes. The inclusion of this analysis permits us to investigate the effect

487 of reference genome origin on both DNA methylation analyses and on marker-based analyses,

488 as opposed to the window-based genome scans reported above. The alignment of reads to

489 the references showed significantly higher efficiency when using a local European reference

490 genome (72.5%) compared to the North American reference (68.4%), resulting in an increase

491 of 6% (paired t-test; t=-44.44, df=49, P<0.001). The more efficient mapping to a local reference

492 produced a significantly higher calling of cytosine bases (6% increase, paired t-test; t=-28.31,

493 df=49, P<0.001) and methylated Cs (8.2% increase, paired t-test; t=-33.18, df=49, P<0.001).

494 Similarly, the comparison of the number of methylated sites per fish after filtering for low

495 coverage sites revealed that using the local European reference genome resulted in identifying

496 more methylated sites (560629.08 ± 12167.89) than with the North America reference

497 (510240.14 ± 11037.78, paired t-test; t=-44.20, df=49, P<0.001). Additionally, the number of

498 Differentially Methylated Sites (DMS) was higher when using a local European reference

499 genome (N=2550 DMS) in comparison to the North American reference (N=2404 DMS). DMS

500 overlapped with 711 and 712 genes in the European and North American assemblies,

501 respectively, and specifically 298 and 299 genes were 1-to-1 orthologs between the two

502 assemblies. A total of 204 of those genes with DMS (68%) were shared in both genomes.

503 There were three significantly enriched GO terms (protein binding, calcium ion binding, and

504 binding) shared in both analyses, and two enriched GO terms (cation binding and metal ion

505 binding) only identified when using a divergent reference genome (Table 2).

# Discussion

Until genome-graphs are more widely available, individual reference genomes remain an integral part of population genomic analyses. However, the reference genome can introduce mapping biases that significantly influence downstream analyses and inferences (Bohling, 2020; Prasad et al., 2022; Valiente-Mullor et al., 2021). Meanwhile, the effects of using a local or more differentiated reference genome remains understudied for ecological and evolutionary model species. To address this knowledge gap, we generated a *de novo* annotated synteny-guided assembly of a European *Gasterosteus aculeatus* fish. Using this novel genome and the established North American reference genome to map sequence data of samples from different populations in Europe and North America, we confirm the reference genome origin significantly impacts downstream analyses. Most notably, a local reference genome increased mapping and genotyping performance. Specifically, mapping efficiency was significantly better using the European reference genome, but the increase in performance was greater for local samples. When using a local reference, more genomic sites were genotyped, and genome window-based estimates of $T_D$ and $F_{ST}$ increased whilst $\pi$ slightly decreased. Similarly, structural variants (SVs) analysis gave slightly different results based on the reference genome used. Consequently, most GO analyses resulted in only a minor proportion of matching enriched GO functions when using different reference genomes. In contrast to the window-based methods, the marker-based DNA methylation analysis pipeline was relatively less affected by reference genome origin, but still about one third of differentially methylated genes was uniquely identified by one but not both references.

## *Genome Assembly of a Gynogenetic Individual*

Recent tools and techniques have increased the efficacy of reference genome assembly, such as utilising long and short read sequencing (Rhie et al., 2021), scaffolding with Hi-C (Peichel et al., 2017), optical or linkage mapping (Glazer et al., 2015), among the growing list of novel and effective techniques (Rhie et al., 2021). The generation of gynogenetic individuals purges genome-wide variation simplifying assembly of genomes (Christensen et al., 2018; Samonte-

533     Padilla et al., 2011). Here, after applying a previously established protocol for gynogenesis,

534     we putatively removed more than 99.9% of genome-wide variation, likely aiding in scaffolding

535     by long-read sequencing. The contiguity of the resulting contig-level gynogen *G. aculeatus*

536     assembly is comparable to the top few assemblies published by the Fish10k project (Fan et

537     al., 2020) in terms of number of contigs and contig N50. Notably, there are a few altered

538     placements of contigs into chromosomes among assemblies. The largest was the 2.07Mb

539     contig tig00002041_pilon, which aligns to both the middle of chrVIII (9.25Mb-10.87Mb) and

540     the end of chrIX (17.89Mb-18.29Mb) and was placed in Gy_chrIX by Chromosemble. The

541     Atlantic and Pacific *G. aculeatus* clades diverged an estimated 44.6Kya (Fang et al., 2018)

542     leaving substantial time for large genomic rearrangements to occur. As such the correct

543     placement of tig00002041_pilon is similarly likely in either chromosome which illustrates the

544     need for further investigations to resolve the differences among assemblies. Overall, the new

545     European gynogen reference genome is high quality, enabling us to test the effects of

546     reference genome origin on downstream population genomic analyses.

### *Mapping and Calling Variants*

548     Despite continuing advances in tools to assemble reference genomes and map sequenced

549     reads, difficulty remains in correctly mapping reads to complex genomic regions enriched in

550     heterozygosity, structural variation, or repetitive elements (Kajitani et al., 2014; Treangen &

551     Salzberg, 2012). By using a reference genome with a longer evolutionary time to the most

552     recent common ancestor (MRCA), complex variants have time to accumulate, likely

553     decreasing mapping efficiency to genomic regions with arguably some of the most sought-

554     after features (e.g., polymorphic regions associated with rapid evolutionary changes and

555     adaptations). Here, we show that using a European stickleback reference genome that has a

556     lower time to the MRCA with European populations increases mapping efficiency and

557     decreases missing data. Conversely, the North American populations also showed increased

558     efficiency when mapped to the European reference, but the difference was noticeably smaller

559     than for European populations. Such a result indicates the removal of heterozygosity through

560  gynogenesis (Samonte-Padilla et al., 2011) and the putative resolution of complex genomic

561  regions in the European assembly improved mapping efficiency. These results are concordant

562  with ethnicity-specific reference genome studies, which have demonstrated that local

563  reference genomes increase depth of coverage resulting in increased sensitivity in variant

564  calling (Ameur et al., 2017; Dewey et al., 2011; Fakhro et al., 2016).

565  *Genome Scans*

566  The effects of improved mapping efficiency and a decrease in missing data were observed to

567  have significant impacts on the estimation of important population genomic metrics. Here,

568  genome-wide estimates of $F_{ST}$ and $T_D$ were higher when using a local reference genome,

569  whereas the opposite pattern was true for estimates of nucleotide diversity (π). Despite the

570  differences in genome-wide estimates, the genome-wide distribution of outlier windows of π,

571  $F_{ST}$, and $T_D$ did not significantly differ across the genome for the same resequenced

572  populations mapped to different reference genomes. Crucially, however, the detection of

573  outlier genes was strongly impacted when using different reference genomes, even when

574  conservatively limiting the analysis to 1-to-1 orthologs identified among the references.

575  Specifically, genes overlapping π outliers had a low proportion of matching orthologs in the

576  same populations when mapped to different reference genomes. Genes overlapping outliers

577  from scans for $F_{ST}$ and $T_D$ were more consistent, but still revealed a large number of

578  differences. In total, only 5.45% of the orthologs among assemblies were mapped to different

579  chromosomes, clearly affecting the detection of outlier genes but only explaining a small

580  proportion of differences observed. The difference in gene placement may instead highlight

581  numerous resolved differences among assemblies that have accumulated since the

582  divergence of Atlantic and Pacific *G. aculeatus* clades 44.6 Kya (Fang et al., 2018). For

583  example, a small deletion in one reference assembly and not the other can result in a gene

584  overlapping an outlier genomic window in only one scan. Overall improved mapping may help

585  specific population genomic analyses including genome-wide representation sequencing

586    where population structure may end up obscuring some SNPs present in only one or few

587    populations (e.g., Baltazar-Soares et al., 2020).

## *Structural Variant Detection*

589    Similar to the effect of reference genome origin on SNP-based scans, the detection of SVs

590    appears to be affected by both reference genome origin and the assembly regardless of origin.

591    Firstly, fewer deletions and inversions were identified when using a local reference genome.

592    This result follows expectations, as there is less time between MRCA for SVs to build up

593    between the sampled population and the reference genome. Secondly, more deletions and

594    duplications and fewer inversions were detected when using the European reference,

595    irrespective of population origin, suggesting differences in the genome assembly plays a role

596    in the detection of SVs. However, the differences in SV detection did not translate into any

597    significant differences in the number of genes overlapping deletions. Overall, our results

598    suggest the detection of SVs may be affected by a combination of mapping efficiency, time to

599    MCRA, and the methods used to assemble a reference genome.

## *Methylation*

601    The most consistent analysis in terms of overlap among reference genomes was the marker-

602    based DNA methylation test. Firstly, using a local reference genome significantly increased

603    mapping efficiency, which resulted in more methylated sites and DMS being detected. The

604    number of genes overlapping the DMS using either reference genomes was the most

605    consistent among all analyses, with one fewer gene (0.14% of total genes) and ortholog

606    (0.34% of total 1-to-1 orthologs) being identified when using a local reference. The DMS

607    analyses recovered a relatively high proportion of overlapping outlier orthologs among

608    reference genomes (68% of 1-to-1 orthologs), but still revealed an effect of the choice of

609    reference. The higher proportion of overlap compared to window-based analyses may be due

610    to the marker-based approach, which is less sensitive to genomic translocation, genome

611    evolution, or assembly errors.

## *GO Enrichment*

The largest effect of reference genome origin was in the GO enrichment analyses of outlier genes from genome scans, with only a minor proportion of enriched GO terms overlapping when using different reference genomes. Given the overall small proportion of overlapping outlier genes, these results were to be expected. GO enrichment analyses are particularly sensitive to minor changes in the number of genes or annotations in lists of genes (Gaudet & Dessimoz, 2017). It should be noted, however, that to allow for a direct comparison of the effects of the different reference genomes, we focused on 1-to-1 orthologs. GO enrichment analyses are common in population genomics (e.g., Chain et al., 2014; Feulner et al., 2015; Reimegård et al., 2017; Liu et al., 2018), and our results highlight reference genome origin strongly impacts such inferences.

## *Conclusions*

Assembling reference genomes is a fast-moving field of research, which sees persistent updates and novel methodologies adopted (Rhie et al., 2021). Hence, variation seen in new reference genomes can reflect both the geographical range of the species distribution but also variation in methodologies used to sequence and assemble the genomes. For example, the two *G. aculeatus* genomes used here originate from samples representing two distinct lineages, the European and North American lineage but also differ significantly in how they were generated. Notably, the North American reference genome (Peichel et al., 2017) is part of a series of updates to the original *G. aculeatus* genome assembly (F. C. Jones et al., 2012). The original assembly used entirely Sanger sequence data (F. C. Jones et al., 2012), compared to the PacBio and Illumina sequence data used for the gynogen genome. As such, only when there is a significant interaction between population and reference origin and no obvious bias in the observations towards either reference can we exclude that the differences in sequencing and assembly methodologies are not the primary cause for the observed patterns.

638  The aim of this study was to investigate the effects of a local reference genome and its effects

639  on downstream analyses. The reference-specific patterns of our results highlight that there is

640  no simple solution. We suggest that the quality of the reference genome and annotations

641  remains the single most important factor when choosing which reference to use. However,

642  when multiple similar quality references are available, a local reference genome offers higher

643  mapping efficiencies and decreases the proportion of missing data. The smallest reference

644  effect among our analyses was for the marker-based methylation analysis, which had

645  markedly more overlap among outliers in comparison to the window-based approach or the

646  SV analysis, but still had over 30% of outliers that did not overlap. Taken together, using a

647  local reference genome should increase the confidence of inferences made within a study,

648  even if the difference is only minor.

## *Data Accessibility and Benefit-Sharing*

650  The raw genomic sequences of the 60 *G. aculeatus* fish were obtained from previous

651  publications (Chain et al., 2014; Feulner et al., 2015), and retrieved from the European

652  Nucleotide Archive, accession number ERP004574. The raw methylome sequences were

653  obtained from a previous publication (Sagonas et al., 2020), retrieved from the NIH genetic

654  sequence database, accession number PRJNA605637.

655  The European-derived gynogen reference genome assembly and their annotations are

656  available from the European Nucleotide Archive (ENA) accession number PRJEB54679.

657  Please note, to comply with ENA submission guidelines the unmapped chromosome (chrUn)

658  was broken into individual contigs, and annotations were lifted over using a custom R script.

659  A total of 316 gene annotations spanned contig boundaries in chrUn and were removed. The

660  original chrUn sequence and annotations are available upon request.

661  Scripts used in the data preparation and analysis pipelines are available online

662  (https://github.com/dthorburn/Origin_Matters).

## *Acknowledgements*

## *Author Contributions*

## *References*

Akalin, A., Franke, V., Vlahoviček, K., Mason, C. E., & Schübeler, D. (2015). Genomation: A toolkit to summarize, annotate and visualize genomic intervals. *Bioinformatics*, *31*(7), 1127–1129. https://doi.org/10.1093/bioinformatics/btu775

Ameur, A., Dahlberg, J., Olason, P., Vezzi, F., Karlsson, R., Martin, M., Viklund, J., Kähäri, A. K., Lundin, P., Che, H., Thutkawkorapin, J., Eisfeldt, J., Lampa, S., Dahlberg, M., Hagberg, J., Jareborg, N., Liljedahl, U., Jonasson, I., Johansson, Å., … Gyllensten, U. (2017). SweGen: A whole-genome data resource of genetic variability in a cross-section of the Swedish population. *European Journal of Human Genetics*, *25*(11), 1253–1260. https://doi.org/10.1038/ejhg.2017.130

Andrews, S. (2010). FASTQC A Quality Control tool for High Throughput Sequence Data. *Babraham Institute*.

Auwera, G. A. van der, Carneiro, M. O., Chris Hartl, R. P., Angel, G. del, Levy-Moonshine, A., Jordan, T., Shakir, K., Roazen, D., Thibault, J., Banks, E., Garimella1, K. v., Altshuler, D., Gabriel, S., & DePristo, M. A. (2014). From FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices pipeline. *Current Protocols in Bioinformatics*.

Baltazar-Soares, M., Klein, J. D., Correia, S. M., Reischig, T., Taxonera, A., Roque, S. M., dos Passos, L., Durão, J., Lomba, J. P., Dinis, H., Cameron, S. J. K., Stiebens, V. A., & Eizaguirre, C. (2020). Distribution of genetic diversity reveals colonization patterns and philopatry of the loggerhead sea turtles across geographic scales. *Scientific Reports*, *10*(1). https://doi.org/10.1038/s41598-020-74141-6

697  Barnett, D. W., Garrison, E. K., Quinlan, A. R., Střmberg, M. P., & Marth, G. T. (2011).
698      Bamtools: A C++ API and toolkit for analyzing and managing BAM files. *Bioinformatics*,
699      *27*(12), 1691–1692. https://doi.org/10.1093/bioinformatics/btr174

700  Bates, D., Mächler, M., Bolker, B. M., & Walker, S. C. (2015). Fitting linear mixed-effects
701      models using lme4. *Journal of Statistical Software*, *67*(1), 1–48.
702      https://doi.org/10.18637/jss.v067.i01

703  Berner, D., Roesti, M., Bilobram, S., Chan, S. K., Kirk, H., Pandoh, P., Taylor, G. A., Zhao,
704      Y., Jones, S. J. M., & Defaveri, J. (2019). De novo sequencing, assembly, and
705      annotation of four threespine stickleback genomes based on microfluidic partitioned
706      DNA libraries. *Genes*, *10*(6), 10–15. https://doi.org/10.3390/genes10060426

707  Bohling, J. (2020). Evaluating the effect of reference genome divergence on the analysis of
708      empirical RADseq datasets. *Ecology and Evolution*, *10*(14), 7585–7601.
709      https://doi.org/10.1002/ece3.6483

710  Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: A flexible trimmer for Illumina
711      sequence data. *Bioinformatics*, *30*(15), 2114–2120.
712      https://doi.org/10.1093/bioinformatics/btu170

713  Chain, F. J. J., Feulner, P. G. D., Panchal, M., Eizaguirre, C., Samonte, I. E., Kalbe, M.,
714      Lenz, T. L., Stoll, M., Bornberg-Bauer, E., Milinski, M., & Reusch, T. B. H. (2014).
715      Extensive Copy-Number Variation of Young Genes across Stickleback Populations.
716      *PLoS Genetics*, *10*(12), 1–18. https://doi.org/10.1371/journal.pgen.1004830

717  Chiang, C., Layer, R. M., Faust, G. G., Lindberg, M. R., Rose, D. B., Garrison, E. P., Marth,
718      G. T., Quinlan, A. R., & Hall, I. M. (2015). SpeedSeq: Ultra-fast personal genome
719      analysis and interpretation. *Nature Methods*, *12*(10), 966–968.
720      https://doi.org/10.1038/nmeth.3505

721  Christensen, K. A., Leong, J. S., Sakhrani, D., Biagi, C. A., Minkley, D. R., Withler, R. E.,
722      Rondeau, E. B., Koop, B. F., & Devlin, R. H. (2018). Chinook salmon (Oncorhynchus
723      tshawytscha) genome and transcriptome. *PLoS ONE, 13*(4), 1–15.
724      https://doi.org/10.1371/journal.pone.0195461

725  Depristo, M. A., Banks, E., Poplin, R., Garimella, K. v., Maguire, J. R., Hartl, C., Philippakis,
726      A. A., del Angel, G., Rivas, M. A., Hanna, M., McKenna, A., Fennell, T. J., Kernytsky, A.
727      M., Sivachenko, A. Y., Cibulskis, K., Gabriel, S. B., Altshuler, D., & Daly, M. J. (2011). A
728      framework for variation discovery and genotyping using next-generation DNA
729      sequencing data. *Nature Genetics*, *43*(5), 491–501. https://doi.org/10.1038/ng.806

730  Dewey, F. E., Chen, R., Cordero, S. P., Ormond, K. E., Caleshu, C., Karczewski, K. J.,
731      Whirl-Carrillo, M., Wheeler, M. T., Dudley, J. T., Byrnes, J. K., Cornejo, O. E., Knowles,
732      J. W., Woon, M., Sangkuhl, K., Gong, L., Thorn, C. F., Hebert, J. M., Capriotti, E.,
733      David, S. P., … Ashley, E. A. (2011). Phased whole-genome genetic risk in a family
734      quartet using a major allele reference sequence. *PLoS Genetics, 7*(9).
735      https://doi.org/10.1371/journal.pgen.1002280

736  Dowle, M., Srinivasan, A., Short, T., Lianoglou, S., Saporta, R., & Antonyan, E. (2015).
737      *data.table: extension of data.frame. R package version 1.9.6*. https://cran.r-
738      project.org/package=data.table

739   Emms, D. M., & Kelly, S. (2019). OrthoFinder: Phylogenetic orthology inference for
740        comparative genomics. *Genome Biology*, *20*(238), 1–14.
741        https://doi.org/10.1101/466201

742   Fakhro, K. A., Staudt, M. R., Ramstetter, M. D., Robay, A., Malek, J. A., Badii, R., Al-Marri,
743        A. A. N., Khalil, C. A., Al-Shakaki, A., Chidiac, O., Stadler, D., Zirie, M., Jayyousi, A.,
744        Salit, J., Mezey, J. G., Crystal, R. G., & Rodriguez-Flores, J. L. (2016). The Qatar
745        genome: A population-specific tool for precision medicine in the Middle East. *Human
746        Genome Variation*, *3*(August 2015), 1–7. https://doi.org/10.1038/hgv.2016.16

747   Fan, G., Song, Y., Yang, L., Huang, X., Zhang, S., Zhang, M., Yang, X., Chang, Y., Zhang,
748        H., Li, Y., Liu, S., Yu, L., Chu, J., Seim, I., Feng, C., Near, T. J., Wing, R. A., Wang, W.,
749        Wang, K., … He, S. (2020). Initial data release and announcement of the 10,000 Fish
750        Genomes Project (Fish10K). *GigaScience*, *9*(8), 1–7.
751        https://doi.org/10.1093/gigascience/giaa080

752   Fang, B., Merilä, J., Ribeiro, F., Alexandre, C. M., & Momigliano, P. (2018). Molecular
753        Phylogenetics and Evolution Worldwide phylogeny of three-spined sticklebacks.
754        *Molecular Phylogenetics and Evolution*, *127*(June), 613–625.
755        https://doi.org/10.1016/j.ympev.2018.06.008

756   Felsenstein, J. (1989). PHYLIP—Phylogeny Inference Package (version 3.2). *Cladistics*, *5*,
757        164–166.

758   Feulner, P. G. D., Chain, F. J. J., Panchal, M., Eizaguirre, C., Kalbe, M., Lenz, T. L., Mundry,
759        M., Samonte, I. E., Stoll, M., Milinski, M., Reusch, T. B. H., & Bornberg-Bauer, E.
760        (2013). Genome-wide patterns of standing genetic variation in a marine population of
761        three-spined sticklebacks. *Molecular Ecology*, *22*(3), 635–649.
762        https://doi.org/10.1111/j.1365-294X.2012.05680.x

763   Feulner, P. G. D., Chain, F. J. J., Panchal, M., Huang, Y., Eizaguirre, C., Kalbe, M., Lenz, T.
764        L., Samonte, I. E., Stoll, M., Bornberg-Bauer, E., Reusch, T. B. H., & Milinski, M.
765        (2015). Genomics of Divergence along a Continuum of Parapatric Population
766        Differentiation. *PLoS Genetics*, *11*(2), 1–18.
767        https://doi.org/10.1371/journal.pgen.1004966

768   Formenti, G., Theissinger, K., Fernandes, C., Bista, I., Bombarely, A., Bleidorn, C., Ciofi, C.,
769        Crottini, A., Godoy, J. A., Höglund, J., Malukiewicz, J., Mouton, A., Oomen, R. A., Paez,
770        S., Palsbøll, P. J., Pampoulie, C., Ruiz-López, M. J., Svardal, H., Theofanopoulou, C.,
771        … Zammit, G. (2022). The era of reference genomes in conservation genomics. *Trends
772        in Ecology and Evolution*, *37*(3), 197–202. https://doi.org/10.1016/j.tree.2021.11.008

773   Frantz, L. A. F., Schraiber, J. G., Madsen, O., Megens, H. J., Cagan, A., Bosse, M., Paudel,
774        Y., Crooijmans, R. P. M. A., Larson, G., & Groenen, M. A. M. (2015). Evidence of long-
775        term gene flow and selection during domestication from analyses of Eurasian wild and
776        domestic pig genomes. *Nature Genetics*, *47*(10), 1141–1148.
777        https://doi.org/10.1038/ng.3394

778   Gan, X., Stegle, O., Behr, J., Steffen, J. G., Drewe, P., Hildebrand, K. L., Lyngsoe, R.,
779        Schultheiss, S. J., Osborne, E. J., Sreedharan, V. T., Kahles, A., Bohnert, R., Jean, G.,
780        Derwent, P., Kersey, P., Belfield, E. J., Harberd, N. P., Kemen, E., Toomajian, C., …
781        Mott, R. (2011). Multiple reference genomes and transcriptomes for Arabidopsis
782        thaliana. *Nature*, *477*(7365), 419–423. https://doi.org/10.1038/nature10414

783    Garrison, E., & Marth, G. (2012). *Haplotype-based variant detection from short-read*
784        *sequencing*. 1–9. http://arxiv.org/abs/1207.3907

785    Gaudet, P., & Dessimoz, C. (2017). *Gene ontology: pitfalls, biases, and remedies* (C.
786        Dessimoz & N. Škunca, Eds.). Humana Press. https://doi.org/10.1007/978-1-4939-
787        3743-1_14

788    Gissi, C., Iannelli, F., & Pesole, G. (2008). Evolution of the mitochondrial genome of
789        Metazoa as exemplified by comparison of congeneric species. *Heredity*, *101*(4), 301–
790        320. https://doi.org/10.1038/hdy.2008.62

791    Glazer, A. M., Killingbeck, E. E., Mitros, T., Rokhsar, D. S., & Miller, C. T. (2015). Genome
792        assembly improvement and mapping convergently evolved skeletal traits in
793        sticklebacks with genotyping-by-sequencing. *G3: Genes, Genomes, Genetics*, *5*(7),
794        1463–1472. https://doi.org/10.1534/g3.115.017905

795    Grabherr, M. G., Russell, P., Meyer, M., Mauceli, E., Alföldi, J., di Palma, F., & Lindblad-Toh,
796        K. (2010). Genome-wide synteny through highly sensitive sequence alignment:
797        Satsuma. *Bioinformatics*. https://doi.org/10.1093/bioinformatics/btq102

798    Grytten, I., Rand, K. D., Nederbragt, A. J., & Sandve, G. K. (2020). Assessing graph-based
799        read mappers against a baseline approach highlights strengths and weaknesses of
800        current methods. *BMC Genomics*, *21*(1), 1–9. https://doi.org/10.1186/s12864-020-
801        6685-y

802    Guan, D., Guan, D., McCarthy, S. A., Wood, J., Howe, K., Wang, Y., Durbin, R., & Durbin, R.
803        (2020). Identifying and removing haplotypic duplication in primary genome assemblies.
804        *Bioinformatics*, *36*(9), 2896–2898. https://doi.org/10.1093/bioinformatics/btaa025

805    Guo, Y., Dai, Y., Yu, H., Zhao, S., Samuels, D. C., & Shyr, Y. (2017). Improvements and
806        impacts of GRCh38 human reference on high throughput sequencing data analysis.
807        *Genomics*, *109*(2), 83–90. https://doi.org/10.1016/j.ygeno.2017.01.005

808    Haenel, Q., Roesti, M., Moser, D., MacColl, A. D. C., & Berner, D. (2019). Predictable
809        genome-wide sorting of standing genetic variation during parallel adaptation to basic
810        versus acidic environments in stickleback fish. *Evolution Letters*, *3*(1), 28–42.
811        https://doi.org/10.1002/evl3.99

812    Hirsch, C. N., Hirsch, C. D., Brohammer, A. B., Bowman, M. J., Soifer, I., Barad, O., Shem-
813        Tov, D., Baruch, K., Lu, F., Hernandez, A. G., Fields, C. J., Wright, C. L., Koehler, K.,
814        Springer, N. M., Buckler, E., Buell, C. R., de Leon, N., Kaeppler, S. M., Childs, K. L., &
815        Mikel, M. A. (2016). Draft assembly of elite inbred line PH207 provides insights into
816        genomic and transcriptome diversity in maize. *Plant Cell*, *28*(11), 2700–2714.
817        https://doi.org/10.1105/tpc.16.00353

818    Holt, C., & Yandell, M. (2011). MAKER2: An annotation pipeline and genome-database
819        management tool for second-generation genome projects. *BMC Bioinformatics*, *12*(1).
820        https://doi.org/10.1186/1471-2105-12-491

821    Huang, Y., Chain, F. J. J., Panchal, M., Eizaguirre, C., Kalbe, M., Lenz, T. L., Samonte, I. E.,
822        Stoll, M., Bornberg-Bauer, E., Reusch, T. B. H., Milinski, M., & Feulner, P. G. D. (2016).
823        Transcriptome profiling of immune tissues reveals habitat-specific gene expression
824        between lake and river sticklebacks. *Molecular Ecology*, *25*(4), 943–958.
825        https://doi.org/10.1111/mec.13520

826    International Human Genome Sequencing Consortium. (2001). Initial sequencing and
827        analysis of the human genome. *Nature*, *409*(6822), 860–921.
828        https://doi.org/10.1038/35057062

829    International Human Genome Sequencing Consortium. (2004). Finishing the euchromatic
830        sequence of the human genome. *Nature*, *431*(7011), 931–945.

831    Jeffares, D. C., Jolly, C., Hoti, M., Speed, D., Shaw, L., Rallis, C., Balloux, F., Dessimoz, C.,
832        Bähler, J., & Sedlazeck, F. J. (2017). Transient structural variations have strong effects
833        on quantitative traits and reproductive isolation in fission yeast. *Nature*
834        *Communications*, *8*(14061), 1–11. https://doi.org/10.1038/ncomms14061

835    Jones, F. C., Grabherr, M. G., Chan, Y. F., Russell, P., Mauceli, E., Johnson, J., Swofford,
836        R., Pirun, M., Zody, M. C., White, S., Birney, E., Searle, S., Schmutz, J., Grimwood, J.,
837        Dickson, M. C., Myers, R. M., Miller, C. T., Summers, B. R., Knecht, A. K., … Kingsley,
838        D. M. (2012). The genomic basis of adaptive evolution in threespine sticklebacks.
839        *Nature*, *484*(7392), 55–61. https://doi.org/10.1038/nature10944

840    Jones, M. R., & Good, J. M. (2016). Targeted capture in evolutionary and ecological
841        genomics. *Molecular Ecology*, *25*(1), 185–202. https://doi.org/10.1111/mec.13304

842    Jones, P., Binns, D., Chang, H. Y., Fraser, M., Li, W., McAnulla, C., McWilliam, H., Maslen,
843        J., Mitchell, A., Nuka, G., Pesseat, S., Quinn, A. F., Sangrador-Vegas, A.,
844        Scheremetjew, M., Yong, S. Y., Lopez, R., & Hunter, S. (2014). InterProScan 5:
845        Genome-scale protein function classification. *Bioinformatics*, *30*(9), 1236–1240.
846        https://doi.org/10.1093/bioinformatics/btu031

847    Kajitani, R., Toshimoto, K., Noguchi, H., Toyoda, A., Ogura, Y., Okuno, M., Yabana, M.,
848        Harada, M., Nagayasu, E., Maruyama, H., Kohara, Y., Fujiyama, A., Hayashi, T., & Itoh,
849        T. (2014). Efficient de novo assembly of highly heterozygous genomes from whole-
850        genome shotgun short reads. *Genome Research*, *24*(8), 1384–1395.
851        https://doi.org/10.1101/gr.170720.113

852    Kao, R. R., Haydon, D. T., Lycett, S. J., & Murcia, P. R. (2014). Supersize me: How whole-
853        genome sequencing and big data are transforming epidemiology. *Trends in*
854        *Microbiology*, *22*(5), 282–291. https://doi.org/10.1016/j.tim.2014.02.011

855    Kim, D., Paggi, J. M., Park, C., Bennett, C., & Salzberg, S. L. (2019). Graph-based genome
856        alignment and genotyping with HISAT2 and HISAT-genotype. *Nature Biotechnology*,
857        *37*(8), 907–915. https://doi.org/10.1038/s41587-019-0201-4

858    King, T., Butcher, S., & Zalewski, L. (2017). *Apocrita - High Performance Computing Cluster*
859        *For Queen Mary University Of London*. https://doi.org/10.5281/ZENODO.438045

860    Koren, S., Walenz, B. P., Berlin, K., Miller, J. R., Bergman, N. H., & Phillippy, A. M. (2017).
861        Canu: Scalable and accurate long-read assembly via adaptive κ-mer weighting and
862        repeat separation. *Genome Research*, *27*(5), 722–736.
863        https://doi.org/10.1101/gr.215087.116

864    Korf, I. (2004). Gene finding in novel genomes. *BMC Bioinformatics*, *5*, 1–9.
865        https://doi.org/10.1186/1471-2105-5-59

866    Krueger, F., & Andrews, S. R. (2011). Bismark: A flexible aligner and methylation caller for
867        Bisulfite-Seq applications. *Bioinformatics*, *27*(11), 1571–1572.
868        https://doi.org/10.1093/bioinformatics/btr167

869    Krzywinski, M., Schein, J., Birol, I., Connors, J., Gascoyne, R., Horsman, D., Jones, S. J., &
870        Marra, M. A. (2009). Circos: An information aesthetic for comparative genomics.
871        *Genome Research*. https://doi.org/10.1101/gr.092759.109

872    Kuznetsova, A., Brockhoff, P. B., & Christensen, R. H. B. (2017). lmerTest Package: Tests in
873        Linear Mixed Effects Models. *Journal of Statistical Software*, *82*(13), 1–26.
874        https://doi.org/10.18637/jss.v082.i13

875    Lacaze, P., Pinese, M., Kaplan, W., Stone, A., Brion, M. J., Woods, R. L., McNamara, M.,
876        McNeil, J. J., Dinger, M. E., & Thomas, D. M. (2019). The Medical Genome Reference
877        Bank: a whole-genome data resource of 4000 healthy elderly individuals. Rationale and
878        cohort design. *European Journal of Human Genetics*, *27*(2), 308–316.
879        https://doi.org/10.1038/s41431-018-0279-z

880    Lai, Y. T., Yeung, C. K. L., Omland, K. E., Pang, E. L., Hao, Y., Liao, B. Y., Cao, H. F.,
881        Zhang, B. W., Yeh, C. F., Hung, C. M., Hung, H. Y., Yang, M. Y., Liang, W., Hsu, Y. C.,
882        Yao, C. te, Dong, L., Lin, K., & Li, S. H. (2019). Standing genetic variation as the
883        predominant source for adaptation of a songbird. *Proceedings of the National Academy
884        of Sciences of the United States of America*. https://doi.org/10.1073/pnas.1813597116

885    Lappalainen, I., Lopez, J., Skipper, L., Hefferon, T., Spalding, J. D., Garner, J., Chen, C.,
886        Maguire, M., Corbett, M., Zhou, G., Paschall, J., Ananiev, V., Flicek, P., & Church, D.
887        M. (2013). DbVar and DGVa: Public archives for genomic structural variation. *Nucleic
888        Acids Research*, *41*, 936–941. https://doi.org/10.1093/nar/gks1213

889    Layer, R. M., Chiang, C., Quinlan, A. R., & Hall, I. M. (2014). LUMPY: A probabilistic
890        framework for structural variant discovery. *Genome Biology*, *15*(6), 1–19.
891        https://doi.org/10.1186/gb-2014-15-6-r84

892    Lenth, R., Singmann, H., Love, J., Buerkner, P., & Herve, M. (2020). emmeans : Estimated
893        Marginal Means, aka Least-Squares Means. In *R package version 1.15-15*.
894        https://doi.org/10.1080/00031305.1980.10483031>.License

895    Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with BWA-
896        MEM. *ArXiv*, 1–3. http://arxiv.org/abs/1303.3997

897    Liu, D., Hunt, M., & Tsai, I. J. (2018). Inferring synteny between genome assemblies: A
898        systematic evaluation. *BMC Bioinformatics*, *19*(1), 1–13.
899        https://doi.org/10.1186/s12859-018-2026-4

900    Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput sequencing
901        reads. *EMBnet.Journal*, *17*(1), 1–3.

902    McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., Garimella,
903        K., Altshuler, D., Gabriel, S., Daly, M., & DePristo, M. A. (2010). The genome analysis
904        toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data.
905        *Genome Research*, *20*(9), 1297–1303. https://doi.org/10.1101/gr.107524.110

906    Nath, S., Shaw, D. E., & White, M. A. (2021). Improved contiguity of the threespine
907        stickleback genome using long-read sequencing. *G3 Genes|Genomes|Genetics*, *0*(0),
908        1–9. https://doi.org/10.1093/g3journal/jkab007

909    Palmer, D. H., & Kronforst, M. R. (2020). A shared genetic basis of mimicry across
910        swallowtail butterflies points to ancestral co-option of doublesex. *Nature
911        Communications*, *11*(1), 1–10. https://doi.org/10.1038/s41467-019-13859-y

912     Parker, J., Tsagkogeorga, G., Cotton, J. A., Liu, Y., Provero, P., Stupka, E., & Rossiter, S. J.
913         (2013). Genome-wide signatures of convergent evolution in echolocating mammals.
914         *Nature*, *502*(7470), 228–231. https://doi.org/10.1038/nature12511

915     Peichel, C. L., Sullivan, S. T., Liachko, I., & White, M. A. (2017). Improvement of the
916         Threespine Stickleback Genome Using a Hi-C-Based Proximity-Guided Assembly.
917         *Journal of Heredity*, *108*(6), 693–700. https://doi.org/10.1093/jhered/esx058

918     Pfeifer, B., Wittelsbürger, U., Ramos-Onsins, S. E., & Lercher, M. J. (2014). PopGenome: An
919         efficient swiss army knife for population genomic analyses in R. *Molecular Biology and
920         Evolution*, *31*(7), 1929–1936. https://doi.org/10.1093/molbev/msu136

921     Pirooznia, M., Goes, F., & Zandi, P. P. (2015). Whole-genome CNV analysis: Advances in
922         computational approaches. *Frontiers in Genetics*, *6*(MAR), 1–9.
923         https://doi.org/10.3389/fgene.2015.00138

924     Prasad, A., Lorenzen, E. D., & Westbury, M. v. (2022). Evaluating the role of reference-
925         genome phylogenetic distance on evolutionary inference. *Molecular Ecology
926         Resources*, *22*(1), 45–55. https://doi.org/10.1111/1755-0998.13457

927     Pritt, J., Chen, N. C., & Langmead, B. (2018). FORGe: Prioritizing variants for graph
928         genomes 06 Biological Sciences 0604 Genetics. *Genome Biology*, *19*(1), 1–16.
929         https://doi.org/10.1186/s13059-018-1595-x

930     Pushkarev, D., Neff, N. F., & Quake, S. R. (2009). Single-molecule sequencing of an
931         individual human genome. *Nature Biotechnology*, *27*(9), 847–850.
932         https://doi.org/10.1038/nbt.1561

933     Quinlan, A. R., & Hall, I. M. (2010). BEDTools: A flexible suite of utilities for comparing
934         genomic features. *Bioinformatics*, *26*(6), 841–842.
935         https://doi.org/10.1093/bioinformatics/btq033

936     R Development Core Team. (2019). *R: A language and environment for statistical
937         computing. R Foundation for Statistical Computing, Vienna, Austria. URL http://www.R-
938         project.org/.* https://doi.org/10.1007/978-3-540-74686-7

939     Raudvere, U., Kolberg, L., Kuzmin, I., Arak, T., Adler, P., Peterson, H., & Vilo, J. (2019).
940         g:Profiler: a web server for functional enrichment analysis and conversions of gene lists
941         (2019 update). *Nucleic Acids Research*, *47*(1), 191–198.
942         https://doi.org/10.1093/nar/gkz369

943     Rausch, T., Zichner, T., Schlattl, A., Stütz, A. M., Benes, V., & Korbel, J. O. (2012). DELLY:
944         Structural variant discovery by integrated paired-end and split-read analysis.
945         *Bioinformatics*, *28*(18), 333–339. https://doi.org/10.1093/bioinformatics/bts378

946     Reid, B. N., Moran, R. L., Kopack, C. J., & Fitzpatrick, S. W. (2021). Rapture-ready darters:
947         Choice of reference genome and genotyping method (whole-genome or sequence
948         capture) influence population genomic inference in Etheostoma. *Molecular Ecology
949         Resources*, *21*(2), 404–420. https://doi.org/10.1111/1755-0998.13275

950     Reid, K., Bell, M. A., & Veeramah, K. R. (2021). *Threespine Stickleback: A Model System
951         For Evolutionary Genomics.* 1–27.

952     Reimegård, J., Kundu, S., Pendle, A., Irish, V. F., Shaw, P., Nakayama, N., Sundström, J.
953         F., & Emanuelsson, O. (2017). Genome-wide identification of physically clustered
954         genes suggests chromatin-level co-regulation in male reproductive development in

955    Arabidopsis thaliana. *Nucleic Acids Research*, *45*(6), 3253–3265.
956    https://doi.org/10.1093/nar/gkx087

957    Rhie, A., McCarthy, S. A., Fedrigo, O., Damas, J., Formenti, G., Koren, S., Uliano-Silva, M.,
958    Chow, W., Fungtammasan, A., Kim, J., Lee, C., Ko, B. J., Chaisson, M., Gedman, G.
959    L., Cantin, L. J., Thibaud-Nissen, F., Haggerty, L., Bista, I., Smith, M., … Jarvis, E. D.
960    (2021). Towards complete and error-free genome assemblies of all vertebrate species.
961    In *Nature* (Vol. 592, Issue April, pp. 737–746).
962    https://doi.org/10.1101/2020.05.22.110833

963    Roach, M. J., Schmidt, S., & Borneman, A. R. (2018). Purge Haplotigs: Synteny Reduction
964    for Third-gen Diploid Genome Assemblies. *BMC Bioinformatics*, *19*(460), 1–10.
965    https://doi.org/10.1101/286252

966    Roesti, M., Kueng, B., Moser, D., & Berner, D. (2015). The genomics of ecological vicariance
967    in threespine stickleback fish. *Nature Communications*, *6*(1), 1–14.
968    https://doi.org/10.1038/ncomms9767

969    Roesti, M., Moser, D., & Berner, D. (2013). Recombination in the threespine stickleback
970    genome - Patterns and consequences. *Molecular Ecology*, *22*(11), 3014–3027.
971    https://doi.org/10.1111/mec.12322

972    Ronco, F., Matschiner, M., Böhne, A., Boila, A., Büscher, H. H., Indermaur, A., el Taher, A.,
973    Malinsky, M., Ricci, V., Kahmen, A., Jentoft, S., & Salzburger, W. (2020). Drivers and
974    dynamics of a massive adaptive radiation in African cichlid fish. *Nature*, 1–6.
975    https://doi.org/10.1038/s41586-020-2930-4

976    Sagonas, K., Meyer, B. S., Kaufmann, J., Lenz, T. L., Häsler, R., & Eizaguirre, C. (2020).
977    Experimental parasite infection causes genome-wide changes in DNA methylation.
978    *Molecular Biology and Evolution*, *37*(8), 2287–2299.
979    https://doi.org/10.1093/molbev/msaa084

980    Samonte-Padilla, I. E., Eizaguirre, C., Scharsack, J. P., Lenz, T. L., & Milinski, M. (2011).
981    Induction of diploid gynogenesis in an evolutionary model organism, the three-spined
982    stickleback (Gasterosteus aculeatus). *BMC Developmental Biology*, *11*, 1–11.
983    https://doi.org/10.1186/1471-213X-11-55

984    Shapiro, M. D., Marks, M. E., Peichel, C. L., Blackman, B. K., Nereng, K. S., Jónsson, B.,
985    Schluter, D., & Kingsley, D. M. (2004). Genetic and developmental basis of evolutionary
986    pelvic reduction in threespine sticklebacks. *Nature*, *428*, 717–723.
987    https://doi.org/10.1038/nature04500

988    Sherry, S. T., Ward, M.-H., Kholodov, M., Baker, J., Phan, L., Smigielski, E. M., & Sirotkin, K.
989    (2001). dbSNP: the NCBI database of genetic variation. *Nucleic Acids Research*, *29*(1),
990    308–311. http://www.ncbi.nlm.nih.gov/SNP.

991    Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. v., & Zdobnov, E. M. (2015).
992    BUSCO: Assessing genome assembly and annotation completeness with single-copy
993    orthologs. *Bioinformatics*, *31*(19), 3210–3212.
994    https://doi.org/10.1093/bioinformatics/btv351

995    Smit, A., Hubley, R., & Green, P. (2015). RepeatMasker Open-4.0. In *RepeatMasker Open-
996    4.0*.

997    Spitz, F., Gonzalez, F., Peichel, C., Vogt, T. F., Duboule, D., & Zákány, J. (2001). Large
998    scale transgenic and cluster deletion analysis of the HoxD complex separate an

999    ancestral regulatory module from evolutionary innovations. *Genes and Development*,
1000    *15*(17), 2209–2214. https://doi.org/10.1101/gad.205701

1001    Springer, N. M., Anderson, S. N., Andorf, C. M., Ahern, K. R., Bai, F., Barad, O., Barbazuk,
1002    W. B., Bass, H. W., Baruch, K., Ben-Zvi, G., Buckler, E. S., Bukowski, R., Campbell, M.
1003    S., Cannon, E. K. S., Chomet, P., Kelly Dawe, R., Davenport, R., Dooner, H. K., Du, L.
1004    H., … Brutnell, T. P. (2018). The maize w22 genome provides a foundation for
1005    functional genomics and transposon biology. *Nature Genetics*, *50*(9).
1006    https://doi.org/10.1038/s41588-018-0158-0

1007    Stamatakis, A. (2014). RAxML version 8: A tool for phylogenetic analysis and post-analysis
1008    of large phylogenies. *Bioinformatics*, *30*(9), 1312–1313.
1009    https://doi.org/10.1093/bioinformatics/btu033

1010    Stanke, M., Diekhans, M., Baertsch, R., & Haussler, D. (2008). Using native and syntenically
1011    mapped cDNA alignments to improve de novo gene finding. *Bioinformatics*, *24*(5), 637–
1012    644. https://doi.org/10.1093/bioinformatics/btn013

1013    Stapley, J., Reger, J., Feulner, P. G. D., Smadja, C., Galindo, J., Ekblom, R., Bennison, C.,
1014    Ball, A. D., Beckerman, A. P., & Slate, J. (2010). Adaptation genomics: The next
1015    generation. *Trends in Ecology and Evolution*, *25*(12), 705–712.
1016    https://doi.org/10.1016/j.tree.2010.09.002

1017    Stern, D. ben, & Lee, C. E. (2020). Evolutionary origins of genomic adaptations in an
1018    invasive copepod. *Nature Ecology and Evolution*, *4*(8), 1084–1094.
1019    https://doi.org/10.1038/s41559-020-1201-y

1020    Treangen, T. J., & Salzberg, S. L. (2012). Repetitive DNA and next-generation sequencing:
1021    Computational challenges and solutions. *Nature Reviews Genetics*, *13*(1), 36–46.
1022    https://doi.org/10.1038/nrg3117

1023    Valiente-Mullor, C., Beamud, B., Ansari, I., Frances-Cuesta, C., Garcia-Gonzalez, N., Mejia,
1024    L., Ruiz-Hueso, P., & Gonzalez-Candelas, F. (2021). One is not enough: On the effects
1025    of reference genome for the mapping and subsequent analyses of short-reads. *PLoS*
1026    *Computational Biology*, *17*(1), 1–29. https://doi.org/10.1371/JOURNAL.PCBI.1008678

1027    Vezzi, F., Narzisi, G., & Mishra, B. (2012). Reevaluating Assembly Evaluations with Feature
1028    Response Curves: GAGE and Assemblathons. *PLoS ONE*, *7*(12), 1–11.
1029    https://doi.org/10.1371/journal.pone.0052210

1030    Walker, B. J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., Cuomo, C. A.,
1031    Zeng, Q., Wortman, J., Young, S. K., & Earl, A. M. (2014). Pilon: An integrated tool for
1032    comprehensive microbial variant detection and genome assembly improvement. *PLoS*
1033    *ONE*, *9*(11). https://doi.org/10.1371/journal.pone.0112963

1034    Wang, Z., Pascual-anaya, J., Zadissa, A., Li, W., Niimura, Y., Huang, Z., Li, C., White, S.,
1035    Xiong, Z., Fang, D., Wang, B., Ming, Y., Chen, Y., Zheng, Y., Kuraku, S., Pignatelli, M.,
1036    Herrero, J., Beal, K., Nozawa, M., … Wang, J. (2013). The draft genomes of soft-shell
1037    turtle and green sea turtle yield insights into the development and evolution of the turtle-
1038    specific body plan. *Nature Genetics*, *45*(6), 701–706. https://doi.org/10.1038/ng.2615

1039    Yu, G., Smith, D. K., Zhu, H., Guan, Y., & Lam, T. T. Y. (2017). Ggtree: an R Package for
1040    Visualization and Annotation of Phylogenetic Trees With Their Covariates and Other
1041    Associated Data. *Methods in Ecology and Evolution*, *8*(1), 28–36.
1042    https://doi.org/10.1111/2041-210X.12628

1043

1044

# Tables and Figures

**Table 1.** Assembly Statistics

| | Contig Assembly | Anchored Assembly |
|---|---|---|
| Number of contigs | 1906 | 877[a] |
| Total size of contigs | 458064828 | 458376166 |
| Longest contig | 5194525 | 33176508 |
| N50 | 746270 | 18858959 |
| **Assembly validation** | | |
| Complete BUSCOs | | 4363 (95.2%) |
| Complete single-copy BUSCOs | | 3913 (85.4%) |
| Complete duplicated BUSCOs | | 450 (9.8%) |
| Fragmented BUSCOs | | 93 (2.0%) |
| Missing BUSCOs | | 128 (2.8%) |
| Total BUSCO groups searched | | 4584 |

[a]21 LGs, chrM, and 855 unplaced contigs.

**Table 2.** Distribution of outlier windows and differentially methylated sites (DMS), overlapping genes, and their functional enrichment.

| Population Origin | Reference Origin | Metric | Outlier Windows | Outlier Genes | 1-1 Ortholog Outlier Genes | Shared Outlier Orthologs | Percentage Overlapping Outliers | Significant GO Terms | Overlapping Terms |
|---|---|---|---|---|---|---|---|---|---|
| North America | North America | Tajima's $D$ | 878 | 1019 | 335 | 221 | 65.97% | 3 | 0 |
| North America | Europe | Tajima's $D$ | 874 | 916 | 333 | | 66.37% | 0 | |
| Europe | North America | Tajima's $D$ | 1308 | 1321 | 342 | 214 | 62.57% | 0 | 0 |
| Europe | Europe | Tajima's $D$ | 1309 | 1449 | 424 | | 50.47% | 3 | |
| North America | North America | $F_{ST}$ | 446 | 541 | 194 | 99 | 51.03% | 0 | NA |
| North America | Europe | $F_{ST}$ | 440 | 534 | 174 | | 56.90% | 0 | |
| Europe | North America | $F_{ST}$ | 672 | 703 | 240 | 101 | 42.08% | 0 | 0 |
| Europe | Europe | $F_{ST}$ | 660 | 788 | 243 | | 41.56% | 2 | |
| North America | North America | π | 888 | 696 | 212 | 29 | 13.68% | 0 | NA |
| North America | Europe | π | 880 | 558 | 153 | | 18.95% | 0 | |
| Europe | North America | π | 1332 | 581 | 143 | 36 | 25.17% | 0 | NA |
| Europe | Europe | π | 1320 | 630 | 199 | | 18.09% | 0 | |
| Europe | North America | DMS | 2404 | 712 | 299 | 204 | 68.23% | 5 | 3 |
| Europe | Europe | DMS | 2550 | 711 | 298 | | 68.45% | 3 | |

**Table 3.** Distribution of structural variants and structurally variable genes and their functional enrichment.

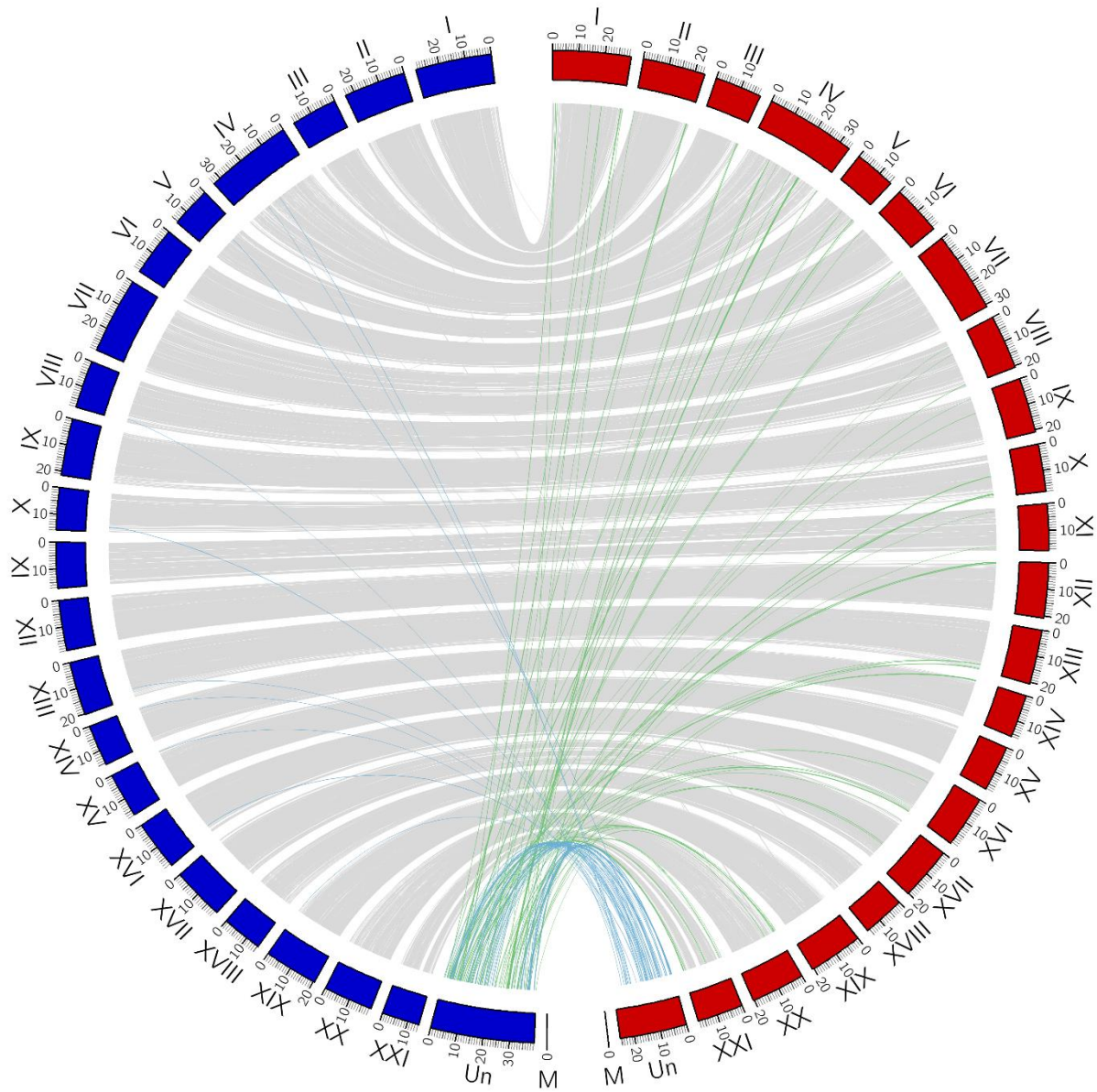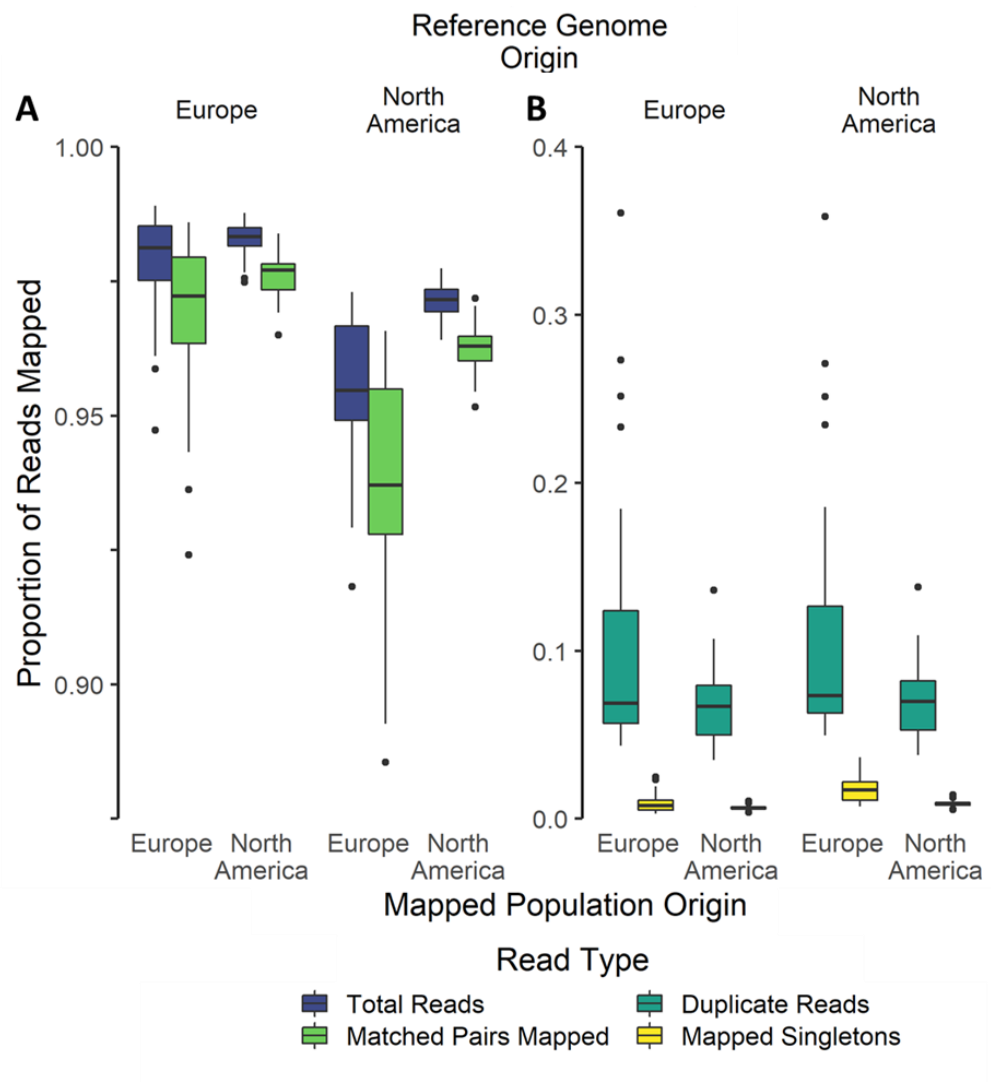| Population Origin | Reference Origin | SV Type | Number of SVs | Number of SVGs | 1-1 Ortholog SVGs | Shared Outlier Orthologs | Percentage Overlapping | Significant GO Terms | Overlapping Terms |
|---|---|---|---|---|---|---|---|---|---|
| North America | North America | Deletion | 3915 | 3034 | 920 | 330 | 35.87% | 3 | 3 |
| North America | Europe | Deletion | 5720 | 2613 | 877 | | 37.63% | 13 | |
| Europe | North America | Deletion | 4049 | 2968 | 893 | 296 | 33.15% | 1 | 0 |
| Europe | Europe | Deletion | 5027 | 2113 | 716 | | 41.34% | 16 | |
| North America | North America | Duplication | 2684 | 3908 | 1195 | 500 | 41.84% | 0 | NA |
| North America | Europe | Duplication | 3402 | 3180 | 1125 | | 44.44% | 0 | |
| Europe | North America | Duplication | 4300 | 4070 | 1256 | 524 | 41.72% | 2 | 0 |
| Europe | Europe | Duplication | 4425 | 3072 | 1074 | | 48.79% | 0 | |
| North America | North America | Inversion | 794 | 4128 | 1358 | 626 | 46.10% | 20 | 0 |
| North America | Europe | Inversion | 757 | 3382 | 1234 | | 50.73% | 10 | |
| Europe | North America | Inversion | 828 | 4031 | 1352 | 522 | 38.61% | 15 | 0 |
| Europe | Europe | Inversion | 685 | 2991 | 1081 | | 48.29% | 0 | |

# Figure Legends

**Figure 1. Synteny plot between the two *G. aculeatus* genome assemblies**. Grey lines between the autosomes represent +99% syntenic blocks greater than 1kb between the European gynogen assembly (left, blue blocks) and the North American *G. aculeatus* assembly (right, red blocks). Coloured lines represent synteny between resolved regions and the unmapped scaffold in each assembly. Specifically, blue and green lines represent 1kb +99% syntenic blocks between the unmapped scaffolds and the alternate reference autosome.

**Figure 2. The effect of reference genome origin on mapping efficiency**. Scales differ among panels, but the units are the same (*A-B*). There was significantly higher mapping efficiency when using the European reference genome regardless of sample origin for total reads (matched and singleton), matched pair reads (*A*) as well as significantly lower singletons (*B*). There was no difference in the number of duplicate reads identified. Reference genome origin is labelled at the top of the plot, and mapped population origin at the bottom.
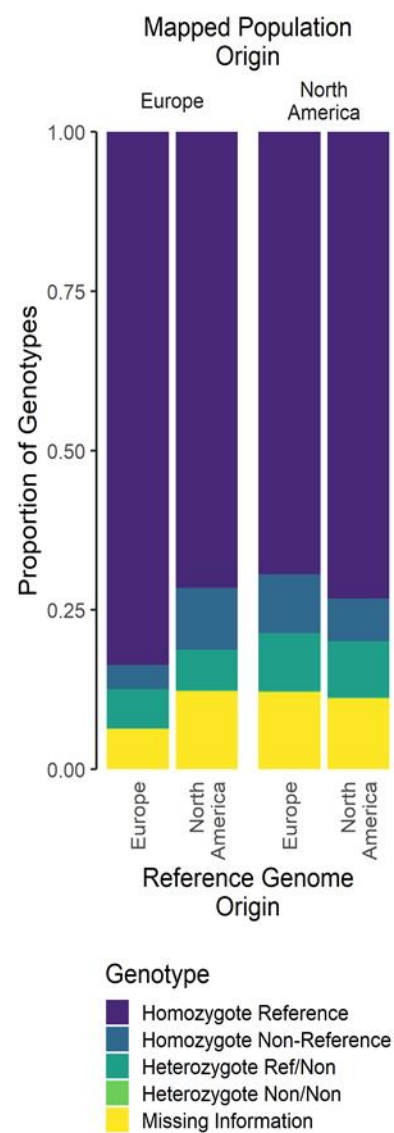
**Figure 3. The proportion of SNP classes among all segregating sites**. Segregating sites with no coverage or with a SNP that was removed during filtering are defined as missing information.
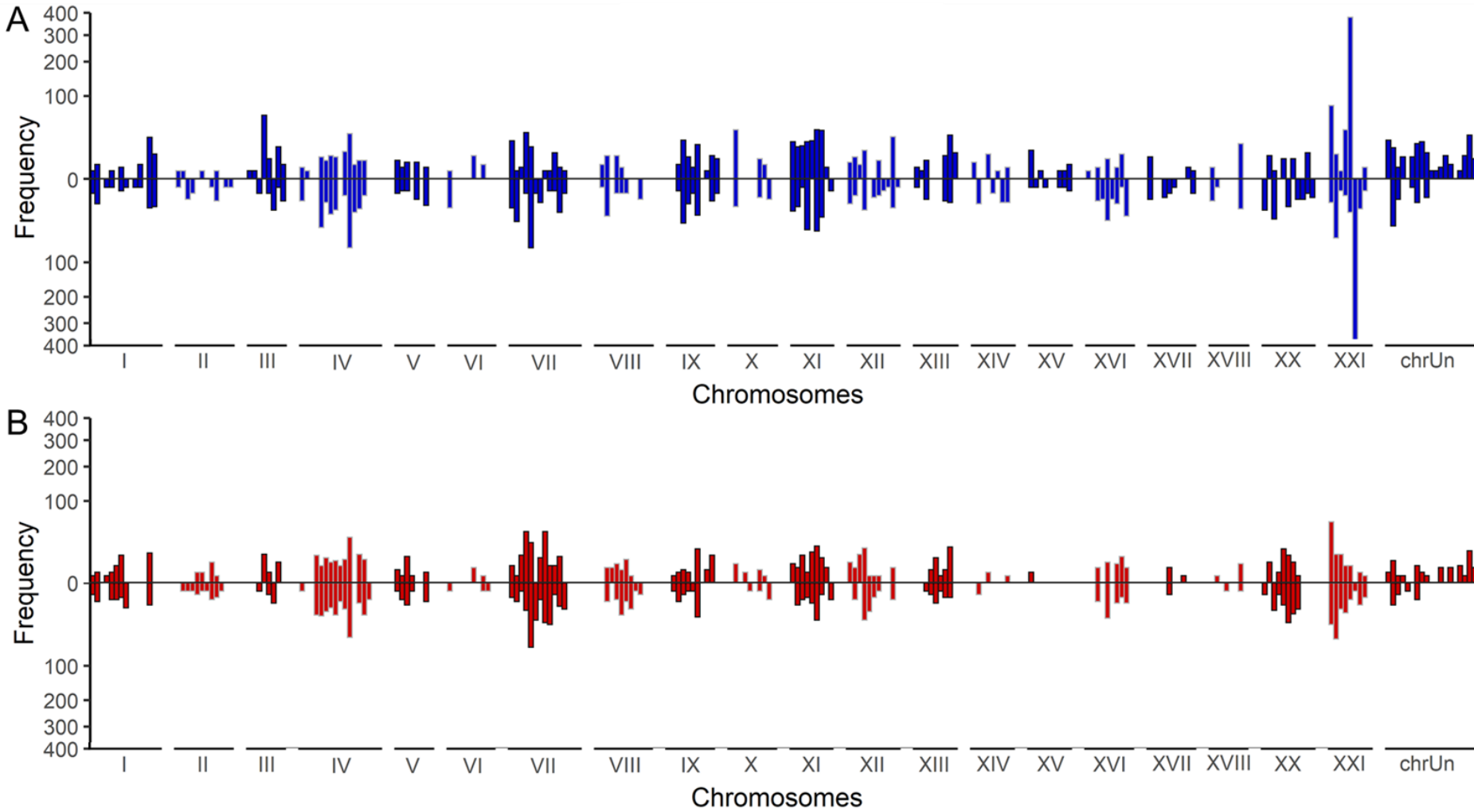
**Figure 4. Comparing distributions of π outliers.** Windows are compared across the genome for (top) European and (bottom) North American populations mapped to the (*A*) European or (*B*) North American reference genome. Axes are square root transformed.

20

21

22

24

25