

1

2 Received Date : 04-Apr-2016

3 Revised Date : 26-May-2016

4 Accepted Date : 31-May-2016

5 Article type : Special Issue

6

7

8 **Corresponding Author Email ID: [eric.archer@noaa.gov](mailto:eric.archer@noaa.gov)**

9 ***strataG*: An R package for manipulating, summarizing, and analyzing**  
10 **population genetic data**

11 Frederick I. Archer<sup>1</sup>, Paula E. Adams<sup>2</sup>, Brita B. Schneiders<sup>3</sup>

12 <sup>1</sup>Southwest Fisheries Science Center, 8901 La Jolla Shores Drive, La Jolla, CA 92037

13 <sup>2</sup>University of Alabama, Box 870344, Tuscaloosa, AL 35487

14 <sup>3</sup>Northwestern University, 2145 Sheridan Road, Evanston, IL 60208

15 **Abstract**

16 We introduce the R package *strataG* as a user-friendly population genetics toolkit.  
17 *strataG* provides easy access to a suite of standard genetic summaries as well as the  
18 ability to rapidly manipulate stratified genetic data for custom analyses. Tests of  
19 population subdivision with most common measures of population subdivision (e.g.,  
20  $F_{ST}$ ,  $G_{ST}$ ,  $\phi_{ST}$ ,  $\chi^2$ ) can be conducted within a single function. The package also  
21 provides wrapper functions that allow users to configure and run popular external

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1111/1755-0998.12559](https://doi.org/10.1111/1755-0998.12559)

This article is protected by copyright. All rights reserved

22 programs such as Genepop, STRUCTURE, and fastsimcoal from within R, and  
23 smoothly interface with popular R packages *adegenet*, and *pegas*. *strataG* is  
24 intended to be an open-source dynamic package that will grow with future needs  
25 and user input.

## 26 **Introduction**

27 The R programming language (R Core Team 2015) has rapidly become a popular  
28 platform for analyses of genetic data. Multiple packages have been developed to  
29 efficiently store and manipulate genetic data (Jombart *et al.* In Prep; Paradis *et al.*  
30 2004) summarize genetic diversity (Goudet 2005; Jombart 2008; Kamvar *et al.*  
31 2014; Keenan *et al.* 2013; Paradis 2010), and conduct phylogenetic analyses  
32 (Paradis *et al.* 2004; Schliep 2011). Still, many population genetics studies that have  
33 relatively standard workflows use a mixed array of software and require the  
34 reformatting of data or results multiple times in order to move among programs.

35 In order to alleviate this movement among software platforms and help users create  
36 custom analytical workflows within the R environment, we present the *strataG*  
37 package, which is designed to be an extensible toolkit for population genetics  
38 analyses. With *strataG*, users can easily compile suites of standard genetic summary  
39 statistics (e.g., allele frequencies, heterozygosity, proportion of unique alleles,  
40 number of private alleles) and conduct common analyses of population structure  
41 (e.g.,  $F_{ST}$ ,  $G_{ST}$ ,  $\phi_{ST}$ ,  $\chi^2$ ) with the flexibility of testing multiple stratification schemes.  
42 The package also includes wrapper functions to run several popular external  
43 analytical programs such as Genepop (Raymond & Rousset 1995) or STRUCTURE  
44 (Pritchard *et al.* 2000). The results of many of these external programs are  
45 automatically read back into the R environment to facilitate downstream analyses  
46 and visualization.

## 47 **Data input and manipulation**

48 Most of the functions in *strataG* operate on genetic data stored within a `gtypes`  
49 object, which is structured as an S4 class with slots for genotypes, stratification

50 schemes, optional sequences, a description of the data, and optional ancillary  
 51 information. A `gtypes` object can be created from data originating from a number of  
 52 standard R data structures, the most common of which is a table-like object (an R  
 53 `matrix` or `data.frame`) organized where each row is an individual and the columns  
 54 are its genotypes. Haplotype IDs from multiple DNA loci can also be stored along  
 55 with their respective sequences.

56 *strataG* has also been designed to work with other popular genetics packages. As  
 57 such, functions (`gtypes2genind`, `gtypes2loci`, and `gtypes2phyDat`) are available  
 58 to convert between `gtypes` and `genind` objects for the *adegenet* package (Jombart  
 59 2008), `loci` objects for the *pegas* package (Paradis 2010), or `phyDat` objects for the  
 60 *phangorn* package (Schliep 2011). Additionally, functions are available to help  
 61 convert and prepare data for loading, such as splitting alleles of a locus that are  
 62 concatenated into a single column, or translating a table of haplotype frequencies to  
 63 a conformable table of haplotype assignment for individuals.

64 Within a `gtypes` object, genotypes are stored as an R `data.frame` of factors which  
 65 makes memory management and calculation of allele frequencies more efficient.  
 66 DNA sequences are also efficiently stored as a `multidna` object from the *apex*  
 67 package (Jombart et al. In Prep), simplifying analyses of multiple haploid loci from  
 68 the same object. A suite of accessor functions is available for `gtypes` objects to  
 69 extract basic information such as the number of individuals, the current  
 70 stratification scheme, or the set of associated sequences. Subsetting or indexing of  
 71 `gtypes` for a specified set of individuals, loci, or strata, uses standard R syntax as in  
 72 this example based on bottlenose dolphin (*Tursiops truncatus*) microsatellite  
 73 genotypes as presented in Lowther-Thieleking *et al.* (2015).

```
74 > # An example microsatellite `gtypes` object
75 > data(msats.g)
76 > msats.g
```

```
77
78 <<< dolphin msats >>>
```

```
79
```

```

80 Contents: 126 samples, 5 loci, 3 strata
81
82 Strata summary:
83     num.samples num.missing num.alleles prop.unique.alleles
84 Coastal        68         1.2         4.8         0.0857
85 Offshore.North  40         0.8        12.6        0.2240
86 Offshore.South  18         0.0        11.0        0.2510
87     heterozygosity
88 Coastal        0.631
89 Offshore.North  0.790
90 Offshore.South  0.867
91
92 Locus summary:
93     num.genotyped num.alleles prop.unique.alleles obsvd.heterozygosity
94 D11t            125         12         0.2500        0.704
95 EV37            119         22         0.1364        0.697
96 EV94            125         15         0.0667        0.776
97 Ttr11           125          9         0.2222        0.704
98 Ttr34           126         10         0.2000        0.698
99 > # Extract the first two loci from the Coastal population
100 > msats.g[, 1:2, "Coastal"]
101
102 <<< dolphin msats >>>
103
104 Contents: 68 samples, 2 loci, 1 stratum
105
106 Strata summary:
107     num.samples num.missing num.alleles prop.unique.alleles
108 Coastal        68          3          5         0.214
109     heterozygosity
110 Coastal        0.571
111
112 Locus summary:
113     num.genotyped num.alleles prop.unique.alleles obsvd.heterozygosity
114 D11t            67          3         0.000        0.522
115 EV37            63          7         0.429        0.619

```

116 A `gtypes` object can also contain alternative stratification schemes for individuals,  
 117 stored as an R `data.frame`. Each column is a unique stratification scheme, with  
 118 individuals in the rows being assigned to a stratum within that scheme. Individuals  
 119 can be excluded from a stratification within a scheme by assigning them the value  
 120 `NA`. Stratification schemes can then be easily changed for different analyses using the  
 121 `stratify` function. For instance, in the following example, the object is restratified  
 122 according to the `broad` scheme which is a stratification column in the `schemes` slot of  
 123 `msats.g`.

```
124 > # Restratify based on "broad" scheme
125 > msats.broad <- stratify(msats.g, "broad")
126 > msats.broad
```

127

```
128 <<< dolphin msats >>>
```

129

```
130 Contents: 126 samples, 5 loci, 2 strata
```

131

```
132 Strata summary:
```

	num.samples	num.missing	num.alleles	prop.unique.alleles
Coastal	68	1.2	4.8	0.0857
Offshore	58	0.8	13.6	0.1751

```
136 heterozygosity
```

Coastal	0.631
Offshore	0.814

139

```
140 Locus summary:
```

	num.genotyped	num.alleles	prop.unique.alleles	obsvd.heterozygosity
D11t	125	12	0.2500	0.704
EV37	119	22	0.1364	0.697
EV94	125	15	0.0667	0.776
Ttr11	125	9	0.2222	0.704
Ttr34	126	10	0.2000	0.698

147 Using standard R functions such as `sapply` or `lapply`, the same analyses can be  
 148 conducted for a suite of stratification schemes with the results for all combined into  
 149 single object for further processing or summary.

## 150 **Summaries**

151 The standard display of a `gtypes` object provides information about the size and  
 152 contents of the object along with many commonly used summary statistics (e.g., the  
 153 number of alleles, observed and expected heterozygosity, allelic richness) for each  
 154 stratum as well as each locus, and differ for genotype and sequence data. If saved to  
 155 an R object, the result of the `summary` function contains information about haplotype  
 156 or allele frequencies. Each of these summary statistics is also available as individual  
 157 functions if only certain ones are desired for a particular application. A pre-defined  
 158 set of by-locus summaries are available from the `summarizeLoci` function.

159 The package also includes summary functions for sequence data, stored either as a  
 160 `gtypes`, `DNABin` (Paradis *et al.* 2004) or `multidna` (Jombart et al. In Prep) object.  
 161 These include functions for calculating transition / transversion ratios ( $T_i/T_v$ ),  
 162 identifying fixed and variable sites, and calculating estimates of selective pressure  
 163 such as Tajima's  $D$  and Fu's  $F_s$  (Tajima 1989; Fu 1997). A set of sequence-specific  
 164 summaries, such as length distributions and base frequencies are available from the  
 165 `summarizeSeqs` function.

```
166 > # Reading a fasta file of aligned mitochondrial control region
167 sequences from Lowther et al (2012) to a DNABin object
168 > fname <- system.file("extdata/dolph.seqs.fasta", package = "strataG")
169 > x <- read.fasta(fname) # one can also use the function
170 read.dna(fname, type = "fasta") in the ape package
171 > x
```

172 126 DNA sequences in binary format stored in a list.

173

174 All sequences of same length: 402

175

```

176 Labels: 4495 4496 4498 5814 5815 5816 ...
177
178 Base composition:
179     a     c     g     t
180 0.301 0.229 0.129 0.341
181
182 > # Summarize sequences
183 > head(summarizeSeqs(x))
184     start end length num.ns num.indels
185 4495    1 402     402         0         2
186 4496    1 402     402         0         2
187 4498    1 402     402         0         1
188 5814    1 402     402         0         2
189 5815    1 402     402         0         2
190 5816    1 402     402         0         2
191
192 > # Calculate transition/transversion ratio
193 > TiTvRatio(x)
194     Ti         Tv Ti.Tv.ratio
195 41.0         4.0         10.2
196
197 > # Estimate Tajima's D test of selective neutrality and p-value
198 > tajimasD(x)
199     D p.value
200 gene.1 -0.506 0.328
201
202 > # For comparison, here is Fu's Fs statistic
203 > fusFs(x)
204 gene.1
205 -7.61
206
207 Quality control checks

```

204 In order to facilitate the use of routine error checks and quality control analyses as  
 205 part of all analytical workflows, we have provided a suite of functions for quality  
 206 assurance / quality control (QA/QC) checks. As an example, there is a function that  
 207 will identify potential duplicate genotypes by reporting all individuals that share  
 208 genotypes across a specified number or fraction of loci (`dupGenotypes`). The  
 209 `jackHWE` function will identify homozygotes that occur at an unusually low  
 210 frequency, implementing the Hardy-Weinberg (HW) jackknife procedure described  
 211 in Morin *et al.* (2010). The result of this function identifies all individuals that, when  
 212 removed from the data, will cause a locus that was previously out of HW equilibrium  
 213 (HWE) to be in HWE. Most by-individual and by-locus summaries and QA/QC checks  
 214 have also been bundled into a single function (`qaqc`) that conducts all tests and will  
 215 optionally write the resulting summaries to comma-delimited (.csv) text files, as  
 216 illustrated in the following example:

```
217 > # Example of checks/summaries of microsatellite data:
218 > checks <- qaqc(msats.g)
219
220 2016-05-25 13:14:26 : Individual summaries
221 2016-05-25 13:14:27 : Locus summaries
222 2016-05-25 13:14:27 : Duplicate genotypes
223 2016-05-25 13:14:30 : Writing files
224 > # By-sample summaries
225 > head(checks$by.sample)
226
227 sample num.loci.missing.genotypes pct.loci.missing.genotypes
228 1 4495 2 0.4
229 2 4496 2 0.4
230 3 4498 0 0.0
231 4 5814 0 0.0
232 5 5815 0 0.0
233 6 5816 0 0.0
pct.loci.homozygous
```



```

234 1          0.667
235 2          0.333
236 3          0.600
237 4          0.200
238 5          0.600
239 6          0.000

240 > # By-locus summaries for Coastal stratum
241 > head(checks$by.locus$Coastal)

242      locus num.genotyped prop.genotyped num.alleles allelic.richness
243 D11t D11t          67          0.985          3          0.0448
244 EV37 EV37          63          0.926          7          0.1111
245 EV94 EV94          68          1.000          5          0.0735
246 Ttr11 Ttr11         68          1.000          4          0.0588
247 Ttr34 Ttr34         68          1.000          5          0.0735

248      prop.unique.alleles exptd.heterozygosity obsvd.heterozygosity
249 D11t          0.000          0.491          0.522
250 EV37          0.429          0.608          0.619
251 EV94          0.000          0.770          0.735
252 Ttr11         0.000          0.662          0.632
253 Ttr34         0.000          0.688          0.647

254 > # Duplicate checks
255 > head(checks$dup.df)

256      ids.1 ids.2 strata.1 strata.2 num.loci.genotyped num.loci.shared
257 1 41579 45237 Coastal Coastal          4          4
258 2 23945 78065 Coastal Coastal          5          4
259 3 25503 78053 Coastal Coastal          5          4
260 4 25509 41822 Coastal Coastal          5          4
261 5 41540 78040 Coastal Coastal          5          4
262 6 41578 45233 Coastal Coastal          5          4

263      prop.loci.shared mismatch.loci
264 1          1.0
265 2          0.8          EV37

```

266	3	0.8	Ttr11
267	4	0.8	Ttr34
268	5	0.8	D11t
269	6	0.8	EV94

## 270 Population Structure

271 At the heart of *strataG* are the tests of population structure. Included in the package  
 272 are functions to calculate  $F_{ST}$ ,  $F'_{ST}$ ,  $G_{ST}$ ,  $G'_{ST}$ ,  $G''_{ST}$ ,  $\theta_{ST}$ ,  $\chi^2$ , and Jost's D. Each function  
 273 takes a `gtypes` object, and returns the test statistic based on the current  
 274 stratification of samples as well as the permutation test p-value, and optionally the  
 275 null distribution of the statistic from the random permutations. Each statistic can be  
 276 run independently or multiple statistics can be run at once with either the  
 277 `overallTest` or `pairwiseTest` functions. The former performs a "global" test of  
 278 population differentiation, while the latter performs tests across all pairs of strata.  
 279 In both `overallTest` and `pairwiseTest`, individual statistics can be specified, or all  
 280 statistics appropriate to the data will be run. The `pairwiseTest` function produces a  
 281 single `data.frame` of all results as well as individual pairwise matrices for each test  
 282 statistic.

```

283 > # Fst test of overall structure
284 > statFst(msats.g, nrep = 1000)
285 $stat.name
286 [1] "Fst"
287
288 $result
289 estimate    p.val
290 0.111807 0.000999
291
292 $null.dist
293 NULL
294
294 > # Pairwise test of four measures
295 > pws <- pairwiseTest(msats.g, stats = c("Fst", "Chi2", "Gst", "D"), nrep =

```

```

296 1000, quietly = TRUE)
297 > print(pws$result)

298                pair.label          strata.1      strata.2
299 1      Coastal (68) v. Offshore.North (40)      Coastal Offshore.North
300 2      Coastal (68) v. Offshore.South (18)      Coastal Offshore.South
301 3 Offshore.North (40) v. Offshore.South (18) Offshore.North Offshore.South
302   n.1 n.2      Fst Fst.p.val  Chi2 Chi2.p.val      Gst Gst.p.val      D
303 1  68  40  0.13064  0.000999 476.8  0.000999  0.0626  0.000999 0.37171
304 2  68  18  0.14641  0.000999 438.9  0.000999  0.0643  0.000999 0.41159
305 3  40  18 -0.00417  0.802198  63.9  0.576424 -0.0126  0.758242 0.00597
306   D.p.val
307 1 0.00105
308 2 0.00108
309 3 0.66152

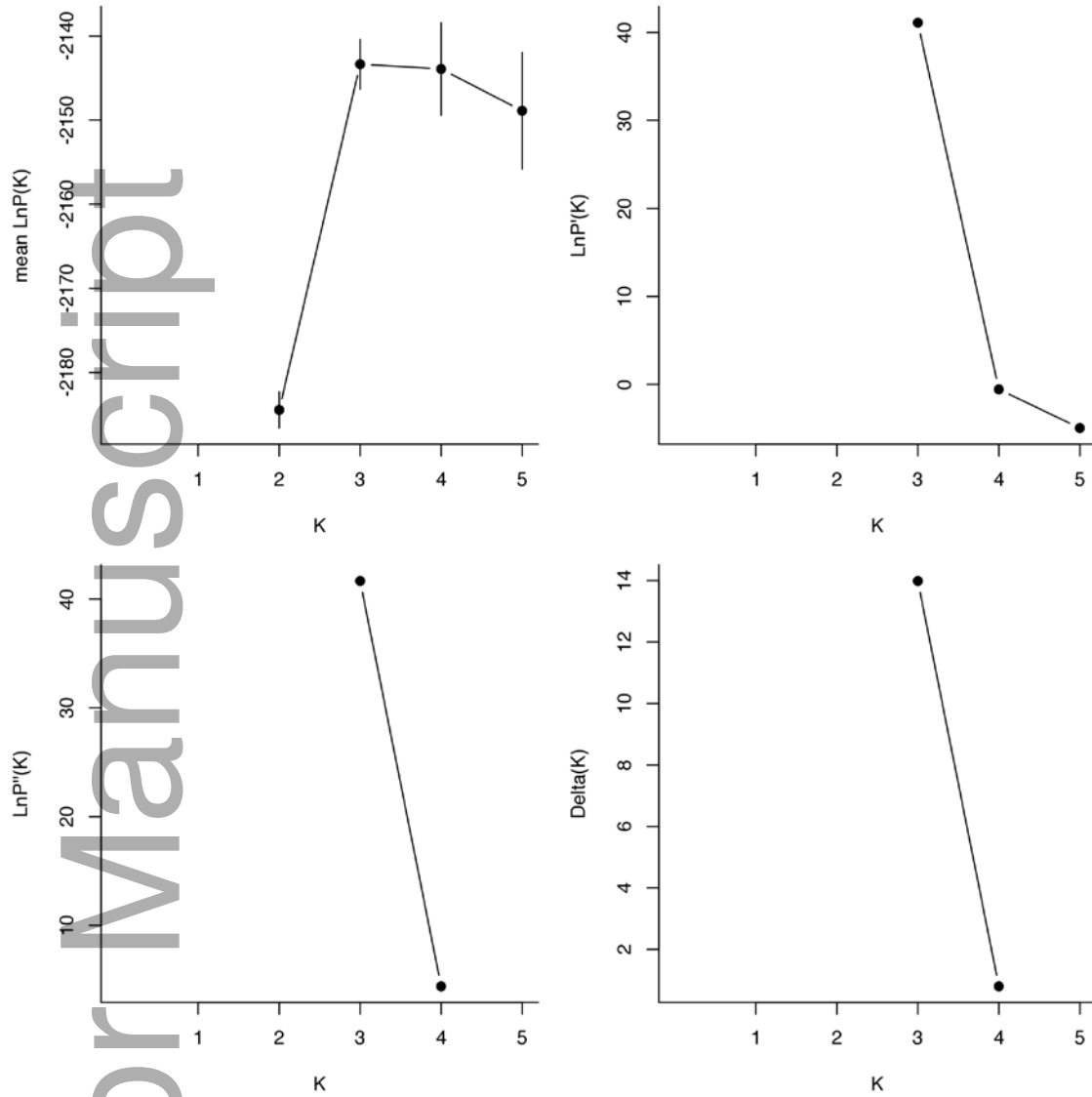
```

### 310 External Software

311 *strataG* provides wrapper functions for several popular population genetics  
312 programs. As of this writing, functions are available for Genepop (Raymond &  
313 Rousset 1995), the Bayesian clustering program STRUCTURE (Pritchard *et al.* 2000),  
314 the DNA sequence alignment program MAFFT (Katoh & Standley 2013), the  
315 coalescent simulator fastsimcoal (Excoffier & Foll 2011), and PHASE (Stephens &  
316 Donnelly 2003) for identifying haplotypes of linked loci. The wrappers are designed  
317 to facilitate use of these programs within the R environment, especially in cases  
318 where R-coded packages with the same functionality are not available. The  
319 wrappers are composed of functions to format and write input data and parameter  
320 files, execute the programs on the command line with options specified as  
321 arguments to the functions, and then in most cases, parse the output from the  
322 programs and return results to the user as R objects. These programs are not  
323 distributed with the *strataG* package and must be downloaded and installed  
324 separately. *strataG* assumes that the programs are installed such that they are  
325 executable on the command line from the working directory. This usually means  
326 installing them into a folder within the system path environmental variable. More  
327 detailed instructions are available in a vignette in the package.

328 Below we demonstrate an example run of STRUCTURE on the bottlenose dolphin  
329 microsatellite data showing differentiation between Coastal and Offshore  
330 populations, but none between Offshore.North and Offshore.South (Lowther-  
331 Thieleking *et al.* 2015). The `structureRun` function formats the microsatellite data,  
332 writes input files, runs the external executable for STRUCTURE, then reads and  
333 parses and reads the output files into an R list structure. The `evanno` function  
334 displays diagnostic plots of number of groups (K) and first and second-order  
335 changes in  $\text{LnP}(K)$  as described in Evanno et al (2005) using base R graphics. The  
336 `c1umpp` function aggregates STRUCTURE runs for a single value of K and is  
337 visualized with the `structurePlot` function which uses the `ggplot2` graphics  
338 package (Wickham 2009).

```
339 > # Run STRUCTURE for k = 2 to 5  
340 > msats.struct <- structureRun(msats.g, k = 2:5, num.k.rep = 100)  
341 >  
342 > # Display diagnostic plots and table from Evanno et al (2005).  
343 > evanno(msats.struct)
```



344

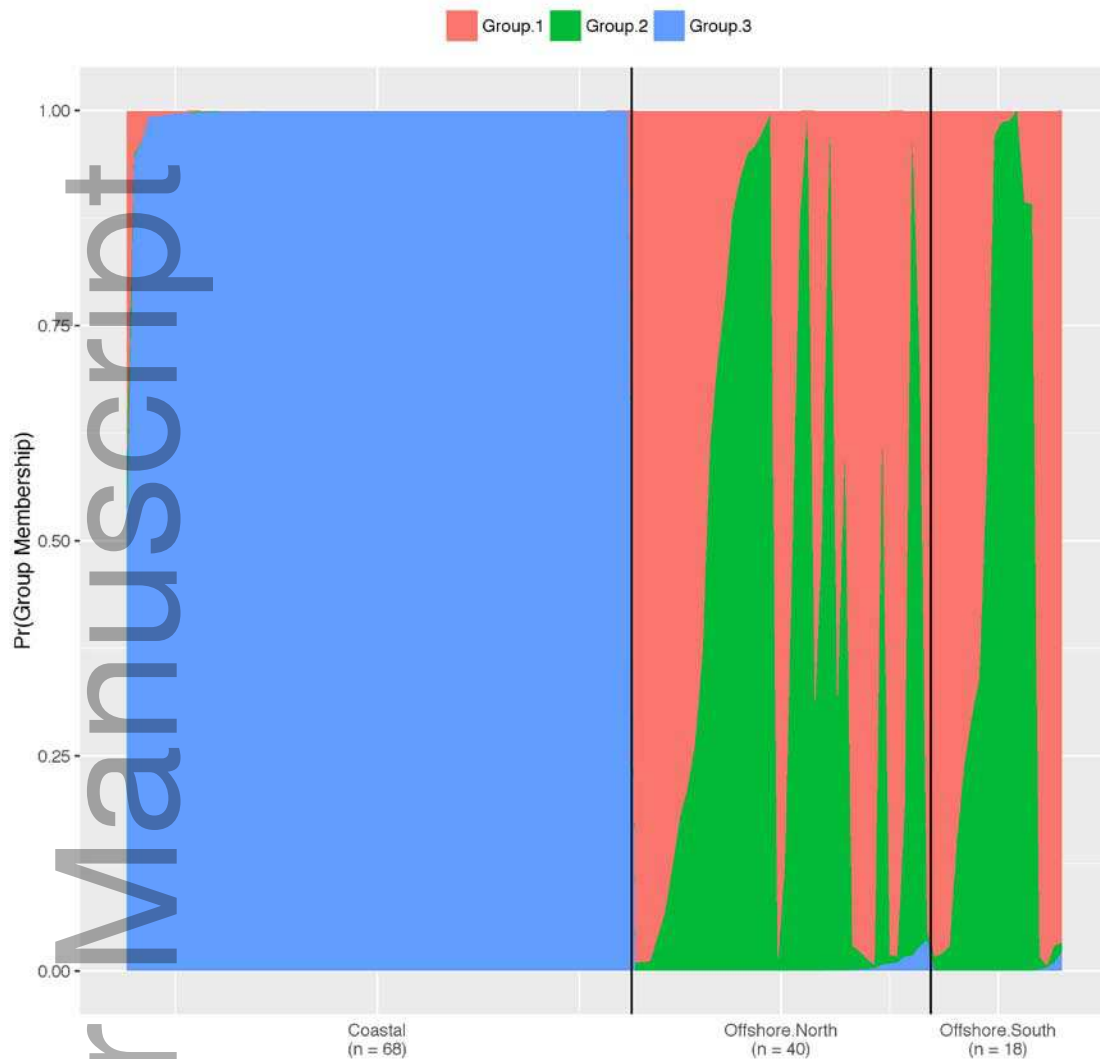
	k	reps	mean.ln.k	sd.ln.k	ln.pk	ln.ppk	delta.k
346	1	2	100	-2184	2.20	NA	NA
347	2	3	100	-2143	3.57	40.545	39.65
348	3	4	100	-2142	5.55	0.898	8.95
349	4	5	100	-2151	6.66	-8.053	NA

350 &gt; # Run CLUMPP to aggregate the results of k = 3

351 &gt; msats.clumpp &lt;- clumpp(msats.struct, k = 3)

352 &gt; # Show distribution of group membership probabilities by strata

353 &gt; sp &lt;- structurePlot(msats.clumpp, horiz = FALSE)



354

355 **Performance**

356 Where possible, *strataG* has been designed for optimal performance in terms of  
 357 computational speed and memory management, while still retaining ease of code  
 358 maintenance. The core algorithms of the computationally-intensive population  
 359 structure tests have been written in C and integrated using the Rcpp package  
 360 (Eddelbuettel 2013). Additionally, where useful, code has been written to take  
 361 advantage of multiple CPUs using functions in the `parallel` package that is  
 362 distributed as part of base R. The functions will automatically detect the user's

363 operating system and set up the clusters in the appropriate format. Users only need  
364 to specify the number of cores to use as a function argument.

365 Although it is difficult to conduct performance benchmark tests that would cover  
366 every dataset and analysis, as a simple example, we generated a simulated  
367 microsatellite dataset of 300 samples, 100 from each of three populations to  
368 demonstrate computational times in *strataG* (Supplemental Materials) Using the  
369 *fastsimcoal* interface in *strataG*, we simulated 1000 loci, then randomly  
370 subsampled for 50, 100, and 500 loci. Mutation rates were randomly chosen for  
371 each loci from a uniform distribution from  $10^{-7}$  to  $10^{-4}$  mutations per generation,  
372 which produced from one to 15 alleles per locus, with a mean of seven. On a  
373 MacBook Pro with a 2.8GHz Intel Core i7 CPU and 16GB of 1600MHz RAM, the  
374 average execution time for calculating overall  $F_{ST}$  and conducting 1000 permutation  
375 replicates to estimate a p-value on each subset as well as all 1000 loci was 0.2s, 0.4s,  
376 2s, and 4s respectively. This indicates a simple linear increase in time with the  
377 number of loci. For most population genetics datasets, processing should be  
378 relatively rapid on standard personal systems or servers.

### 379 **Obtaining *strataG***

380 The current stable release of *strataG* (version 1.0.5 as of this writing) is available for  
381 download from the Comprehensive R Archive Network (CRAN) at [[https://cran.r-](https://cran.r-project.org)  
382 [project.org](https://cran.r-project.org)]. Pre-release versions with recent bug fixes and additions are available  
383 through GitHub at [<https://github.com/EricArcher/strataG>], which also has  
384 instructions on how to install the package using the `install_github` function  
385 available in the `devtools` package.

386 *strataG* is actively maintained, and contains many other additional functions. Users  
387 are encouraged to explore the package starting with the list provided by  
388 `help(package = "strataG")`. The package includes vignettes covering how to  
389 create `gtypes`, summarize data, conduct population differentiation tests, and install  
390 and run external programs. With the rapid spread of next generation sequencing  
391 and the growth of population genomics, we expect an increasing demand to store

392 and process ever larger datasets and the development of novel analytical methods.  
393 Because the `gtypes` object is a strongly-typed S4 object with accessor functions for  
394 the data, the underlying structures can easily be modified to take advantage of more  
395 efficient storage methods without necessitating changes to existing code. Given the  
396 open source nature of the R programming environment and the growing popularity  
397 of collaborative development platforms such as GitHub, we envision *strataG* to be a  
398 dynamic toolkit that can grow with the addition of new methods and community  
399 needs. We invite the population genetics community to actively participate in its  
400 development with suggestions for improvements and contributed code.

#### 401 **Acknowledgements**

402 The authors wish to thank the members of the Marine Mammal Genetics Group at  
403 the Southwest Fisheries Science Center and the many other users who have helped  
404 in the development of *strataG* through their feedback on early versions of the  
405 package. We would also like to thank the National Evolutionary Synthesis Center  
406 (NESCent) for organizing the Population Genetics in R Hackathon, which was held in  
407 March 2015 at the National Evolutionary Synthesis Center (NESCent) in Durham,  
408 NC, with the goal of addressing interoperability, scalability, and workflow building  
409 challenges for the population genetics package ecosystem in R. FIA was a participant  
410 in the hackathon, and is indebted to NESCent (NSF #EF-0905606) for hosting and  
411 supporting the event.

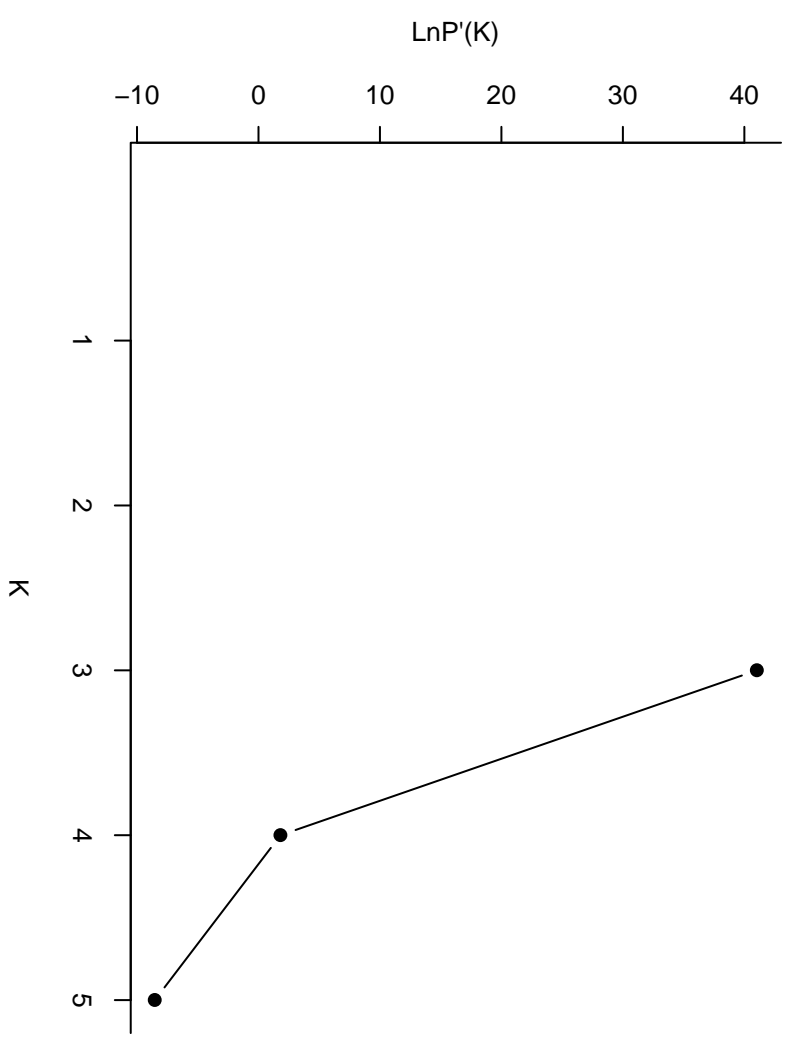
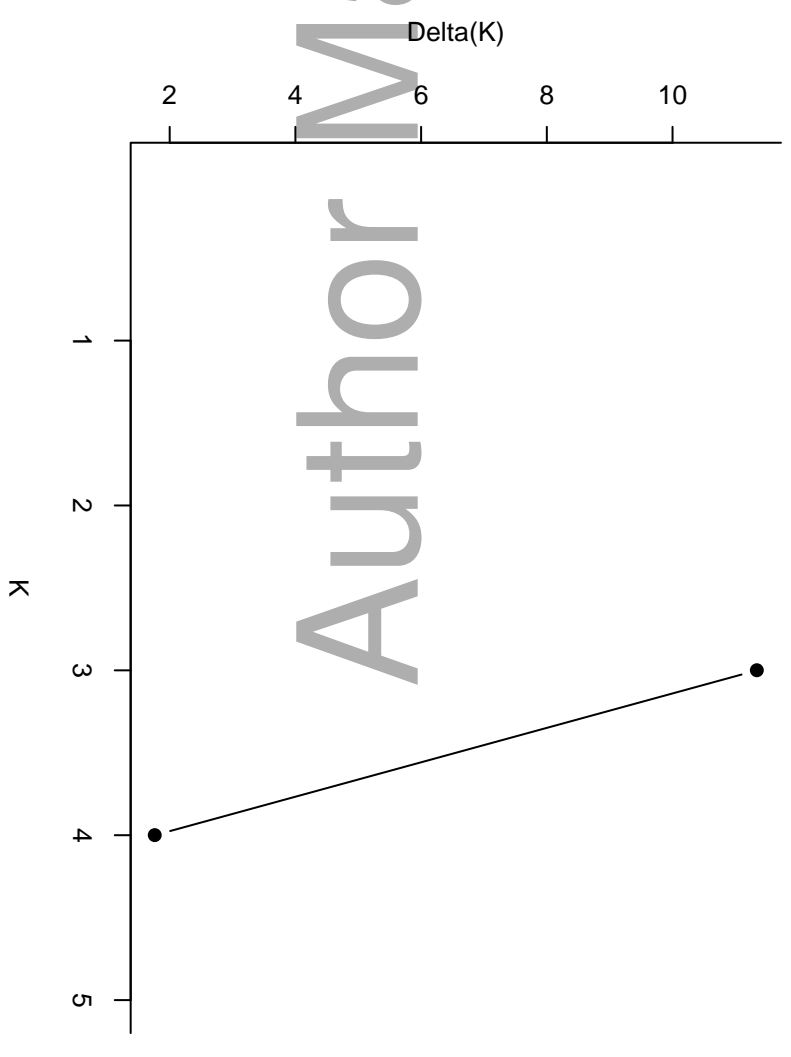
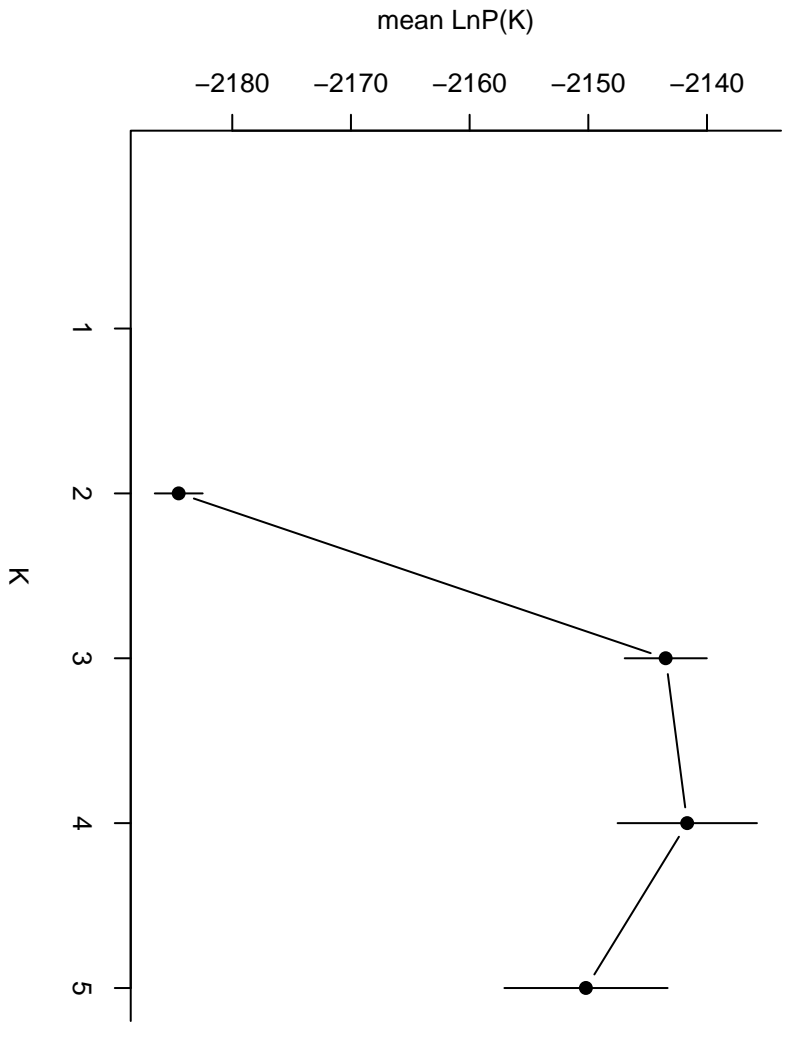
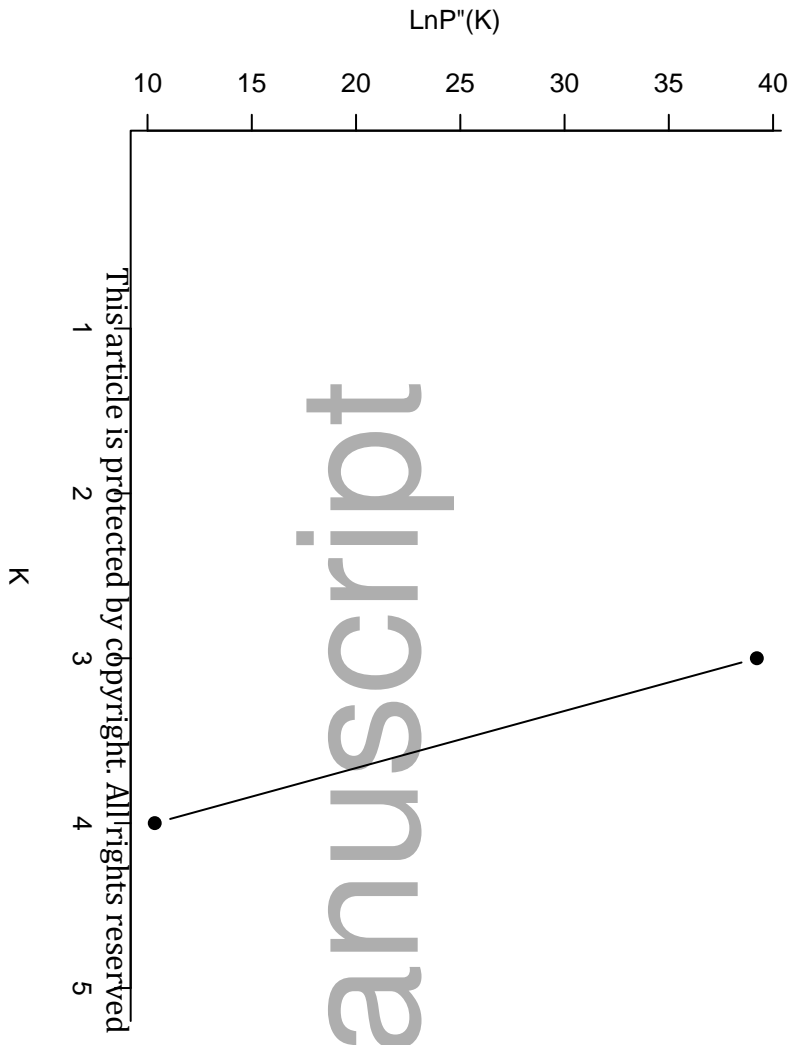
#### 412 **Literature Cited**

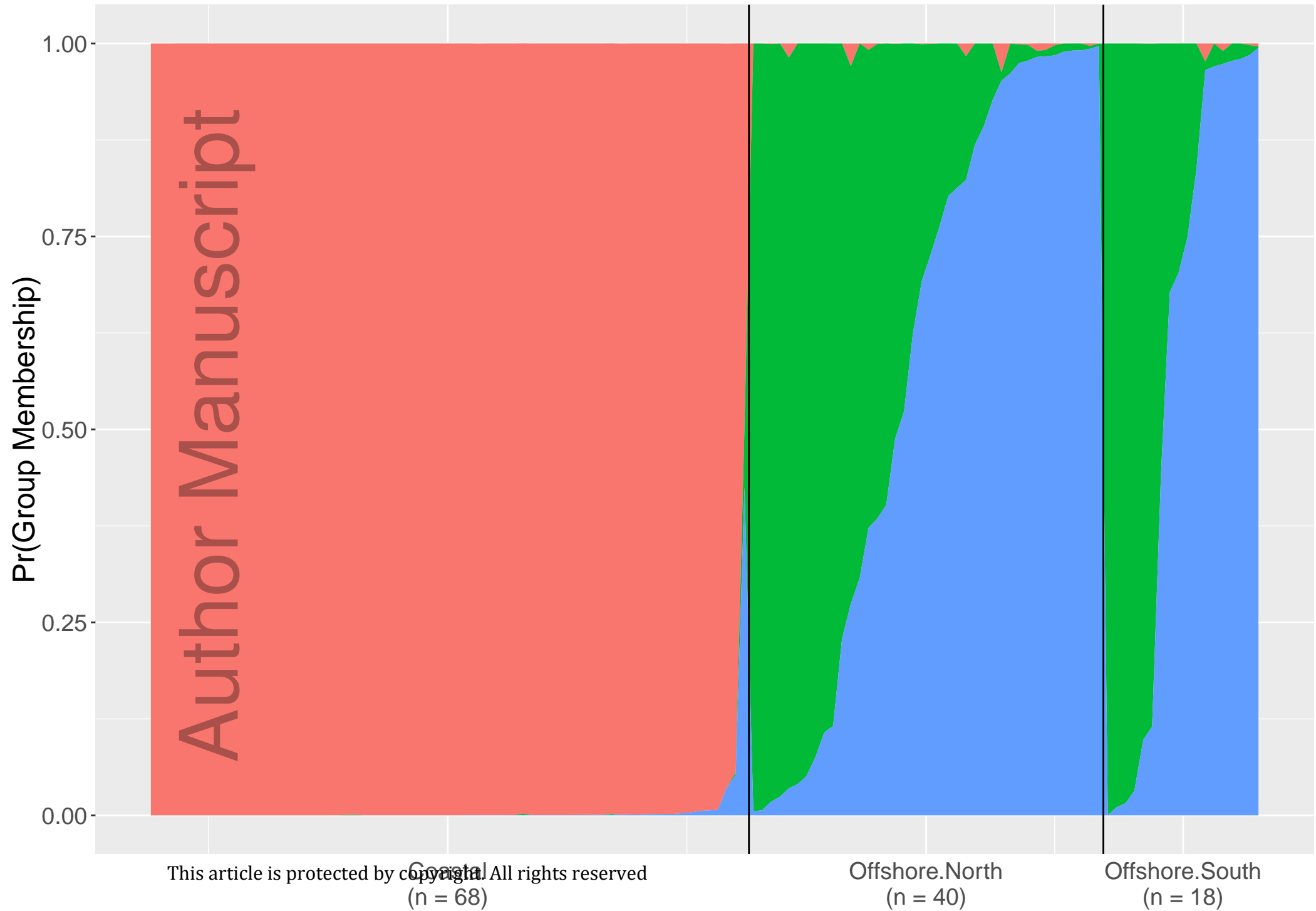
- 413 Eddebuettel D (2013) *Seamless R and C++ Integration with Rcpp*. Springer.
- 414 Excoffier L, Foll M (2011) Fastsimcoal: a continuous-time coalescent simulator of  
415 genomic diversity under arbitrarily complex evolutionary scenarios.  
416 *Bioinformatics* **27**, 1332-1334.
- 417 Fu Y-X (1997) Statistical tests of neutrality of mutations against population growth,  
418 hitchhiking and background selection. *Genetics* **147**,915-925.



- 419 Goudet J (2005) HIERFSTAT, a package for R to compute and test hierarchical F-  
420 statistics. *Molecular Ecology Notes* **5**, 184-186.
- 421 Jombart T (2008) adegenet: a R package for the multivariate analysis of genetic  
422 markers. *Bioinformatics* **24**, 1403-1405.
- 423 Jombart T, Schliep KP, Archer FI, *et al.* (In Prep) apex: phylogenetics with multiple  
424 genets. *Molecular Ecology Resources*.
- 425 Kamvar ZN, Tabima JF, Grünwald NJ (2014) *Poppr*: an R package for genetic analysis  
426 of populations with clonal, partially clonal, and/or sexual reproduction. *PeerJ*  
427 **2**, e281.
- 428 Katoh K, Standley DM (2013) MAFFT Multiple Sequence Alignment Software  
429 Version 7: Improvements in Performance and Usability. *Molecular Biology*  
430 *and Evolution* **30**, 772-780.
- 431 Keenan K, McGinnity P, Cross TF, Crozier WW, Prodöhl P (2013) diveRsity: An R  
432 package for the estimation and exploration of population genetics  
433 parameters and their associated errors. *Methods in Ecology and Evoluton* **4**,  
434 782-788.
- 435 Lowther-Thieleking JL, Archer FI, Lang AR, Weller DW (2015) Genetic  
436 differentiation among coastal and offshore common bottlenose dolphins,  
437 *Tursiops truncatus*, in the eastern North Pacific Ocean. *Marine Mammal*  
438 *Science* **31**, 1-20.
- 439 Morin PA, Archer FI, Foote AD, *et al.* (2010) Complete mitochondrial genome  
440 phylogeographic analysis of killer whales (*Orcinus orca*) indicates multiple  
441 species. *Genome Research* **20**, 908-916.
- 442 Paradis E (2010) pegas: an R package for population genetics with an integrated-  
443 modular approach. *Bioinformatics* **26**, 419-420.
- 444 Paradis E, Claude J, Strimmer K (2004) APE: Analyses of Phylogenetics and  
445 Evolution in R language. *Bioinformatics* **20**, 289-290.

- 446 Pritchard JK, Stephens M, Donnelly P (2000) Inference of population structure using  
447 multilocus genotype data. *Genetics* **155**, 945-959.
- 448 R Core Team (2015) R: A Language and Environment for Statistical Computing. R  
449 Foundation for Statistical Computing.
- 450 Raymond M, Rousset F (1995) GENEPOP (Version 1.2): Population genetics  
451 software for exact tests and ecumenicism. *Heredity* **86**, 248-249.
- 452 Schliep KP (2011) phangorn: phylogenetic analysis in R. *Bioinformatics* **27**, 592-593.
- 453 Stephens M, Donnelly P (2003) A comparison of Bayesian methods for haplotype  
454 reconstruction from population genotype data. *American Journal of Human*  
455 *Genetics* **73**, 1162-1169.
- 456 Tajima F (1989) Statistical method for testing the neutral mutation hypothesis by  
457 DNA polymorphism. *Genetics* **123**, 585-595.
- 458 Wickham H (2009) *ggplot2: Elegant Graphics for Data Analysis*. Springer.





This article is protected by copyright. All rights reserved  
Onsite (n = 68)

Offshore.North (n = 40)

Offshore.South (n = 18)