



**QUEEN'S  
UNIVERSITY  
BELFAST**

## **DiveRsity: An R package for the estimation and exploration of population genetics parameters and their associated errors**

Keenan, K., Mcginnity, P., Cross, T. F., Crozier, W. W., & Prodöhl, P. A. (2013). DiveRsity: An R package for the estimation and exploration of population genetics parameters and their associated errors. *Methods in Ecology and Evolution*, 4(8), 782-788. <https://doi.org/10.1111/2041-210X.12067>

**Published in:**  
Methods in Ecology and Evolution

**Document Version:**  
Peer reviewed version

**Queen's University Belfast - Research Portal:**  
[Link to publication record in Queen's University Belfast Research Portal](#)

**Publisher rights**  
© 2018 British Ecological Society. This work is made available online in accordance with the publisher's policies. Please refer to any applicable terms of use of the publisher.

**General rights**  
Copyright for the publications made accessible via the Queen's University Belfast Research Portal is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**  
The Research Portal is Queen's institutional repository that provides access to Queen's research output. Every effort has been made to ensure that content in the Research Portal does not infringe any person's rights, or applicable UK laws. If you discover content in the Research Portal that you believe breaches copyright or violates any law, please contact [openaccess@qub.ac.uk](mailto:openaccess@qub.ac.uk).

**Open Access**  
This research has been made openly available by Queen's academics and its Open Research team. We would love to hear how access to this research benefits you. – Share your feedback with us: <http://go.qub.ac.uk/oa-feedback>

1 **diveRsity: An R package for the estimation and exploration of**  
2 **population genetics parameters and their associated errors.**

3  
4 Kevin Keenan<sup>1</sup>, Philip McGinnity<sup>2</sup>, Tom F. Cross<sup>2</sup>, Walter W. Crozier<sup>3</sup>, Paulo A. Prodöhl<sup>1,\*</sup>

5  
6 <sup>1</sup>Institute for Global Food Security, School of Biological Science, Medical Biology  
7 Centre, Queen's University, Belfast, BT9 7BL Northern Ireland

8 <sup>2</sup>Aquaculture & Fisheries Development Centre, School of Biological, Earth &  
9 Environmental Sciences, University College Cork, Ireland

10 <sup>3</sup>Agri-Food and Biosciences Institute, Newforge Lane, Belfast, Northern Ireland

11  
12 \*Correspondence author: Paulo A. Prodöhl

13 Address: Institute for Global Food Security, School of Biological Science, Medical  
14 Biology Centre, Queen's University, Belfast, BT9 7BL Northern Ireland

15 E-mail: p.prodohl@qub.ac.uk

16  
17 Running title: diveRsity package

18 Word count: 3102  
19  
20  
21  
22

23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45

## Summary

1. We present a new R package, `diveRsity`, for the calculation of various diversity statistics, including common diversity partitioning statistics ( $\theta$ ,  $G_{ST}$ ) and population differentiation statistics ( $D_{Jost}$ ,  $G'_{ST}$ ,  $\chi^2$  test for population heterogeneity), among others. The package calculates these estimators along with their respective bootstrapped confidence intervals for loci, sample population pairwise and global levels. Various plotting tools are also provided for a visual evaluation of estimated values, allowing users to critically assess the validity and significance of statistical tests from a biological perspective.
2. `diveRsity` has a set of unique features, which facilitate the use of an informed framework for assessing the validity of the use of traditional F-statistics for the inference of demography, with reference to specific marker types, particularly focusing on highly polymorphic microsatellite loci. However, the package can be readily used for other codominant marker types (e.g. allozymes, SNPs).
3. A detailed example of usage and descriptions of package capabilities are provided. The example demonstrates useful strategies for the exploration of data and interpretation of results generated by `diveRsity`. Additional on-line resources for the package are also described, including a GUI web app version intended for those with more limited experience using R for statistical analysis.

## 46 Introduction

47

48 As a consequence of the growing suite of statistical genetics tools, which are often tailored to  
49 particular marker types, the analyses of population genetic data is becoming an increasingly  
50 complex task (Excoffier & Heckel, 2006). For instance, F-statistics is a commonly used  
51 framework for the description of genetic diversity partitioning within and among populations.  
52 F-statistics estimators (e.g.  $\theta$ ,  $G_{ST}$ ) suffer from an incompatibility when applied to highly  
53 polymorphic microsatellite markers (Hedrick, 1999; Jost, 2008), as a result of their negative  
54 dependence on within sub-population heterozygosity (Jost, 2008). Thus, for loci with many  
55 alleles (e.g.  $>10$ ), within sub-population heterozygosity will invariably be high, and as a  
56 consequence, “traditional” F-statistics will have a theoretical maximum well below the  
57 expected  $F_{ST} = 1$ . Attempts have been made to overcome this issue, most notably by  
58 Hedrick (2005), with the development of  $G'_{ST}$  and more recently Jost (2008) with the  
59 development of  $D_{Jost}$ . However, much confusion still exists about what these “new” statistics  
60 should actually be used for (Gerlach *et al.*, 2010). It is not the purpose of this study to  
61 elaborate on such issues, however, interested readers are encouraged to see Jost (2008),  
62 Meirmans & Hedrick (2011) and Whitlock (2011) for useful reviews.

63

64 To add to the complexity, recent advances in molecular screening methodologies have greatly  
65 facilitated the ease with which genetic data can be generated. As a consequence, an  
66 increasing number of researchers, often with a limited background in statistical genetics  
67 analyses (Karl *et al.*, 2012), face the difficult task of analysing and interpreting such data. Thus,  
68 software tools that facilitate this task, by providing suitable frameworks to allow for informed  
69 analysis pipelines are essential. To this end, we present the software `diveRsity`. This R

70 package allows the estimation of various population genetic summary statistics including the  
71 two “traditional” F-statistics analogues;  $\theta$  (Weir & Cockerham, 1984) and  $G_{ST}$  (Nei &  
72 Chesser, 1983), and the two “new” differentiation statistics;  $G'_{ST}$  (Hedrick, 2005) and  $D_{Jost}$   
73 (Jost, 2008), as well as their unbiased/nearly unbiased estimators. Each statistic can be  
74 estimated for locus, global and sample pairwise comparisons. The package also provides  
75 functionality for the estimation of 95% confidence intervals at all relevant levels, through an  
76 integrated bootstrapping procedure. Uniquely to `diveRcity`, various plotting functions,  
77 designed to allow researchers to assess the validity of using their particular data set (or suite  
78 of marker loci) for the inference of geneflow using the F-statistics framework, are also  
79 provided, as well as visualisation tools for large pairwise matrices of genetic differentiation  
80 and parameter confidence intervals. Furthermore, `diveRcity` also provides a range of other  
81 statistical tools, which are commonly used in population genetic analyses pipelines but are  
82 rarely integrated into a single software package.

83

84 Another major advantage of using `diveRcity` is that it produces summary data structures,  
85 which are very close to publication-ready formats (e.g. figure 1). Given that the compilation  
86 of such summary data is time consuming and often involves the use of several software  
87 packages, `diveRcity` offers a valuable addition to the molecular ecologist’s statistical  
88 toolkit. Its implementation as an R package also makes `diveRcity` ideal for easy  
89 incorporation into analysis pipelines where batch processing of files/data is required, as is  
90 often the case in simulation based studies.

91

92 This package is intended to promote a more considered and simplified approach to  
93 frequentist population genetic structure analyses. Through the inclusion of `diversity`

94 partitioning statistics (e.g.  $\theta$  &  $G_{ST}$ ), differentiation statistics (e.g.  $G'_{ST}$  &  $D_{Jost}$ ), as well as  
95 functionality to assess the behaviour of these statistics across loci and population samples,  
96 we hope to give researchers the necessary tools to make educated decisions about the  
97 statistical and biological validity of their analyses with relative ease. Following this rationale,  
98 we have also opted to omit the option for users to carry out  $p$ -value null hypothesis testing in  
99 relation to F-statistics and population sample differentiation estimators. This decision was  
100 taken given the lack of meaningful information conveyed through the use of  $p$ -values in this  
101 context, as well as the many misconceptions that exist regarding the biological interpretation  
102 of  $p$ -values in relation to these statistics (Wagenmakers, 2007). We have instead provided  
103 functions to allow users to estimate 95% confidence intervals (calculated as the 2.5% and  
104 97.5% quantiles of a bootstrap distribution), for a range of statistical estimators calculated by  
105 the package, thus, leading to more reliable conclusions about the biological significance of  
106 trends in the data, (see figure 2 in du Prel *et al.*, 2009), leaving less room for erroneous  
107 interpretation.

108

## 109 **Description**

110

111 `diveRsity` is a package written for use in R (R Development Core Team, 2011). It is primarily  
112 designed for the estimation, exploration and validation of genetic differentiation/structure  
113 indices. The package aims to consolidate under the same work environment, many of the  
114 most popular population genetic statistics such as those mentioned above, in order to provide  
115 researchers with a simplified way in which to calculate and compare these statistics. This  
116 strategy is particularly useful for the identification of polymorphism based biases mentioned  
117 previously. This information can be subsequently used, along with additional exploration tools

118 implemented in the package, to make informed decisions about which statistical measures or  
119 molecular markers can be appropriately applied to address a particular question.

120

121 `diveRcity` also calculates a plethora of other statistics and has various other population  
122 genetics applications. Table 1 provides a list of functions along with brief descriptions of their  
123 specific purposes. The package accepts raw genotype data for any group of co-dominant  
124 molecular markers in the *genepop* file format (Raymond & Rousset, 1995). There is no limit  
125 to the size of the accepted input file other than the amount of random access memory (RAM)  
126 available to users. In addition to providing users with the ability to efficiently estimate an  
127 array of population genetic statistics, `diveRcity` is also particularly flexible in terms of  
128 return result formats (e.g. text files, excel workbooks and native R objects such as matrices  
129 and data frames). This flexibility facilitates subsequent downstream analysis (e.g.  
130 incorporation into simulation or Approximate Bayesian Computation (ABC) pipelines as the  
131 summary statistic calculation software). A list of specific output formats is also summarised  
132 in Table 1.

133

#### 134 **Dependencies and suggested packages**

135

136 In general, `diveRcity` can be used with a standard R installation and two additional  
137 extension packages (`plotrix` and `shiny`). The functions `divPart`, `inCalc`, `chiCalc` and  
138 `readGenepop`, `divBasic`, `bigDivPart` and `divRatio`, (i.e. the major analytical  
139 functions), can all operate independently of non-standard packages. The only disadvantages  
140 of this approach are slower execution times (i.e. parallel computation is not available), and a  
141 limited number of formats available for returned results. To fully capitalise on the additional

142 features of `diveRcity` (listed in Table 1), the installation of all suggested packages is  
143 recommended. Details of these packages are given in Table 2.

144

#### 145 **Comparisons with other software**

146

147 The main motivation behind the development of `diveRcity` was to provide a cross-platform  
148 software, which allows comprehensive and fast frequentist analysis of co-dominant molecular  
149 data, while maintaining usability and convenient result formats. On each of these aims,  
150 `diveRcity` performs comparatively better in relation to other similar software.

151

152

#### 153 ***Comprehensiveness***

154 When compared to other software which estimate similar statistics, `diveRcity` generally  
155 provides a more comprehensive range of parameter calculation options. In terms of the total  
156 number of available population genetics statistics, with the possible exception of the Mac OS  
157 X only program, `GenoDive` (Meirmans & Van Tienderen, 2004), `diveRcity` estimates many  
158 more than `DEMEtics` (Gerlach *et al.*, 2010), `SMOGD` (Crawford, 2010), `mmod` (Winter, 2012),  
159 `hierfstat` (Goudet, 2004) or `SPADE` (Chao & Shen, 2003).

160

161 Focusing only on diversity partitioning/differentiation statistics, `diveRcity` overlaps in its  
162 calculation of  $D_{Jost}$  with all of the above mentioned software. However, `diveRcity` is the  
163 only package that allows the estimation of 95% confidence intervals, globally (i.e. for all  
164 samples and loci), per locus (i.e. over all samples) and for all pairwise sample comparisons  
165 (i.e. over all loci per population pair). `SMOGD`, for example, which is perhaps the most popular



166 of these applications (with over 212 citations according to Google scholar), calculates  
167 bootstrapped confidence intervals for  $D_{Jost}$  at the locus level across all population samples,  
168 but does not provide this estimation for either the global or pairwise levels.

169

170 Despite the focus of this study on diversity partition/differentiation statistics, `diveRcity`  
171 also estimates many other useful population genetics statistics. These include,  $\chi^2$  tests of  
172 Hardy-Weinberg equilibrium (HWE), Allelic richness ( $A_r$ ), Chi-square tests for sample  
173 homogeneity, 'Yardstick' diversity standardised ratios (Skrbinšek *et al.*, 2012) and locus  
174 informativeness for the inference of ancestry (Rosenberg *et al.*, 2003). Contrary to other  
175 similar programs, `diveRcity` also provides various exploratory plotting tools, which can be  
176 very useful for the identification of meaningful trends within results with minimal effort (e.g.  
177 Example 1). Typically, this task would involve the compilation of output results from various  
178 programs and subsequent visualisation in an independent software package (e.g. Microsoft  
179 Excel). A full description of `diveRcity`'s functionality can be found by typing either of the  
180 following commands into the R console:

181

```
182 # diveRcity must be installed
```

183

```
184 # 1) package help pages
```

```
185 help(package = "diveRcity")
```

186

```
187 # 2) package user manual
```

```
188 vignette("diveRcity")
```

189

190

191 **Speed**

192

193 Given the different analytical focuses of distinct softwares, performance comparisons in  
194 terms of speed are not straightforward. For example, while in one software a given test  
195 statistic might be estimated using a maximum likelihood procedure, in another, a more  
196 computational intensive procedure (e.g. bootstrapping) may be used. For the purposes of this  
197 study, comparisons were restricted to instances where distinct softwares implemented similar  
198 computational processes to calculate a similar suit of statistical parameters. Based on these  
199 criteria, only two truly comparable speed comparisons were possible between `diveRcity`  
200 and any of the above listed software.

201

202 The first is a comparison of locus confidence interval estimation using bootstrapping with  
203 SMOGD. The reproducible code used to run `diveRcity` is:

204

```
205 system.time({  
206 # load diveRcity  
207 library("diveRcity")  
208 # load Test_data  
209 data(Test_data)  
210  
211 # run the analysis  
212 x <- divPart(infile = Test_data, outfile = NULL, gp = 3,  
213             pairwise = TRUE, WC_Fst = FALSE, bs_locus = TRUE,  
214             bs_pairwise = FALSE, bootstraps = 1000, plot = FALSE,  
215             parallel = TRUE)  
216 })
```

217 When running SMOGD on the example data set `Test_data` (see Keenan *et al.*, in press for  
218 details on these data), with bootstraps set to 1000, the time taken to return results to the  
219 web browser is 2 min 34.1 sec, while `diveRsity` takes only 1 min 17.3 sec to carry out the  
220 same calculations on a laptop with an Intel Core i5-2435 CPU @ 2.49GHz. It is also relevant to  
221 note that `diveRsity`'s performance can be significantly increased with the use of additional  
222 CPUs.

223

224 The second comparison involves the calculation of diversity partitioning statistics per locus  
225 for large data sets (e.g. RAD-seq derived SNP genotypes). This comparison was carried out  
226 between the `diveRsity` function `bigDivPart` and the `hierfstat` function  
227 `basic.stats`. For this test, a simulated data set of 268 individuals across four population  
228 samples genotyped for 55,200 bi-allelic SNP loci was used. To complete the entire analysis,  
229 `diveRsity` took 3 min 20.1 sec, while `hierfstat` took 6 min 44.8 sec, using the same laptop  
230 as described above. Such speed differences become even more important with the increasing  
231 rate at which large arrays of loci can be genotyped for large numbers of individuals.

232

### 233 ***Usability & convenience***

234

235 Similar to other R packages, in order to fully benefit from all features built into `diveRsity`,  
236 a reasonable level of expertise in R is required. However, `diveRsity` has been designed so  
237 that even R beginners or those with very limited expertise, can easily carry out comprehensive  
238 analysis of their data, including results being written to file, in many cases with a single  
239 command line. This is in contrast to other packages such as `mmod` and `hierfstat` which  
240 invariably require users to export their own result from the R environment, as well as execute

241 more functions to calculate fewer parameters than `diveRsity`. An example of the  
242 convenient results formats returned by `diveRsity` is shown in figure 1.

243

244 In keeping with the focus on ease of use, `diveRsity` also includes a web application which  
245 provides a browser based user interface for the estimation of the most popular statistics  
246 implemented in the command line version of the package. This application was built using the  
247 framework provided by the R package, `shiny` (RStudio & Inc., 2012) and provides users with  
248 a range of benefits including an easy to use interface and downloadable result files. The  
249 browser user interface also allows users to run their analyses on a remote server, thus, local  
250 system resources are not consumed. The application can be accessed at:

251

252 <http://glimmer.rstudio.com/kkeenan/diveRsity-online/>

253

254 Users can also run this application locally by executing the following command in the R  
255 console:

256

```
257 # after loading diveRsity
```

```
258 divOnline()
```

259

260

261 Despite an emphasis on simplicity, `diveRsity` still retains all of the functionality and  
262 flexibility provided by the R environment (i.e. all results are returned to the current session  
263 workspace). Thus, users with more experience, can easily pipe results from their analyses into  
264 downstream custom analyses (e.g. ABC).

265

## 266 **Accessing the package**

267

268 The `diveRcity` package is hosted on the Comprehensive R Archive Network (CRAN), and can  
269 be downloaded using the `install.packages` function in R. Simply type the following  
270 command into the R console:

271

```
272 install.packages("diveRcity", dependencies = TRUE)
```

273

274 Providing the user has a working internet connection, and following the selection of a suitable  
275 CRAN repository mirror, the package will download and install automatically.

276

277 Ongoing development of `diveRcity` can also be tracked at:

278 <http://diversityinlife.weebly.com/software.html>

279 This web page contains the latest developmental versions of the package as well as an update  
280 log.

281

## 282 **Examples**

283

284 As a demonstration of some of the envisaged applications of `diveRcity`, two reproducible  
285 examples are provided below. These examples assume that the `diveRcity`, `shiny`,  
286 `doParallel`, `sendplot` and `plotrix` packages have been installed as well as their  
287 dependencies. For additional examples, users are encouraged to read the package manual.

288

289 **Example 1. Using visualisation tools to investigate large genetic differentiation matrices**

290

291 Pairwise genetic differentiation is an important parameter in the assessment of relationships  
292 among populations within a geographical context. To date, the true potential of pairwise  
293 genetic differentiation statistics has not been fully realised, owing mainly to difficulties in  
294 identifying meaningful trends in often very large numbers of population comparisons.

295

296 However, by using both the `divPart` and `difPlot` functions, `diveRcity` allows users to  
297 visualise large pairwise matrices of genetic differentiation, making the identification of  
298 particularly differentiated population samples relatively straightforward. This procedure is  
299 demonstrated below.

300

301 Load `diveRcity` into the current R session:

302

```
303 # Load the diveRcity package  
304 require("diveRcity")
```

305

306 In this example the `Big_data` data set (distributed with `diveRcity`), will be used. The data  
307 were simulated under a hierarchical island model (i.e. five island groups with 10 sub-  
308 populations each allowing high geneflow within island groups and low geneflow among island  
309 groups), using the software EASYPOP v1.7 (Balloux, 2001). Population samples within the  
310 `Big_data` data file were arranged in order of geographical proximity for the purpose of  
311 demonstrating how `diveRcity` can be used to identify broad-scale geographical trends from  
312 genetic data.

313

```
314 # Load 'Big_data'  
315 data(Big_data, package = "diveRsim")
```

316

317 The `divPart` function is first used to calculate the required pairwise statistics matrices. In  
318 this example the argument `parallel` will be set to `TRUE` as a large number of comparisons  
319 have to be computed (i.e.  $[\frac{1}{2}N] \times [N - 1] = 1225$  for  $N = 50$ ).

320

```
321 # Assign the results to the variable 'pwStats'  
322 # (i.e. pw = pairwise)  
323 pwStats <- divPart(infile = Big_data, outfile = "Big_results",  
324                   gp = 2, WC_Fst = TRUE, bs_locus = FALSE,  
325                   bs_pairwise = FALSE, bootstraps = 0,  
326                   Plot = FALSE, parallel = TRUE)
```

327

328 The resulting R object, `pwStats` contains the required pairwise statistics which can be passed  
329 to the function `difPlot` for visualisation.

330

```
331 difPlot(x = pwStats, outfile = "Big_results",  
332         interactive = TRUE)
```

333

334 This command will write four `.png` files (one for each estimated statistic), and four `.html` files  
335 to the folder `Big_results` under the current R working directory. An example of the  
336 functionality of the `.html` tool-tips is given in figure 2. From this figure, it is clear that the data  
337 are represented by five distinct genetic groups, which correlates with the simulation

338 conditions described above. There are clearly high levels of differentiation among island  
339 groups (light blue/white) and low levels of differentiation within island groups (dark blue).  
340 This graphical representation perfectly relays what is known to be genetically/evolutionarily  
341 true (though natural population systems will rarely be so ideal).  
342 Figure 2 also illustrates the ability to rapidly identify population pairs of interest by simply  
343 positioning the mouse pointer over a particular comparison square/pixel. In this example the  
344 pairwise comparison between populations 18 vs 23, ( $G_{ST} = 0.8883$ ,  $\theta = 0.9408$ ,  $G'_{ST} =$   
345  $0.9927$  and  $D_{Jost} = 0.8802$ ), indicates that these two populations are highly differentiated  
346 from one another.

347

#### 348 **Example 2. Assessing polymorphism bias in diversity partitioning estimators**

349

350 As discussed above, diversity partitioning statistics such as  $G_{ST}$  and  $\theta$  are negatively  
351 dependent on within sub-population heterozygosity. Where this negative dependence is  
352 present (e.g. when using highly polymorphic microsatellites), it is important to ensure that  
353 inferences made from calculated values do not violate important assumptions. Using the  
354 functions `divPart`, `readGenepop` and `corPlot`, it is possible to carry out an *ad hoc*  
355 assessment of polymorphism bias in diversity statistics, thus allowing users to make informed  
356 decisions about whether to proceed with inference of demographic processes for example. A  
357 reproducible example is given below:

358

```
359 # Load the diveRsity package  
360 require("diveRsity")
```

361



362 Next an example data set (`Test_data`) provided with `diveRsimy` should be loaded into the  
363 R session.

364

```
365 # Load 'Test_data'  
366 data(Test_data, package = "diveRsimy")
```

367

368 Initially `Test_data` is analysed by the function `divPart` to calculate locus  $\theta$ ,  $G_{ST}$ ,  $G'_{ST}$   
369 and  $D_{Jost}$  estimators.

370

```
371 # Assign the results to the variable 'difStats'  
372 difStats <- divPart(infile = Test_data, outfile = "Test",  
373                   gp = 3, WC_Fst = TRUE, bs_locus = TRUE,  
374                   bs_pairwise = FALSE, bootstraps = 1000,  
375                   plot = TRUE, parallel = TRUE)
```

376

377 Next `Test_data` is analysed by `readGenepop` to count the total number of alleles per locus.

378

```
379 # Assign the result to the variable 'numAlleles'  
380 numAlleles <- readGenepop(infile = Test_data, gp = 3,  
381                          bootstrap = FALSE)
```

382

383 The package has now generated two results objects in the R environment: `difStats` and  
384 `numAlleles`. These objects can be passed to the function `corPlot`.

385

```
386 corPlot(x = numAlleles, y = difStats)
```

387

388 Figure 3 provides an example of the output from this analysis. As can be seen in this example,  
389 both  $\theta$  and  $G_{ST}$  are negatively correlated with the number of alleles per locus, whilst  $G'_{ST}$   
390 and  $D_{Jost}$  are strongly positively correlated. This discordance is indicative of a case where  
391 the mutation rate is likely to obscure past demographic processes (e.g. geneflow), thus such  
392 a data set is unsuitable for addressing such questions.

393

394 Users executing the above code will also see a range of other graphical outputs in a folder  
395 named "Test" within their working directory. These plots allows users to assess the  
396 variability of parameter estimation for individual loci, which can in turn be incorporated into  
397 decisions about 'misbehaving' loci for example.

398

### 399 **Acknowledgements**

400 The authors would like to thank J.J. Magee, M.S.P Ravinet, J. Coughlan and C. Johnston for  
401 testing the `diveRsim` package and R. Hynes for proofreading the manuscript. We would  
402 also like to express our gratitude to MEE executive editor Dr. Robert B. O'Hara and two  
403 anonymous reviewers, whose comments greatly improved the manuscript and the  
404 `diveRsim` package. K.K. was supported by a PhD studentship from the Beaufort Marine  
405 Research Award in Fish Population Genetics funded by the Irish Government under the Sea  
406 Change programme. P.A.P, T.F.C, W.W.C and P.McG were also supported by this award.

407

### 408 **References**

409

410 Balloux, F. ( 2001) EASYPOP (version 1.7): a computer program for population genetics  
411 simulations. *Journal of Heredity*, **92**, 301–302.

412

413 Chao, A. & Shen, T.J. (2003) Program SPADE (species prediction and diversity estimation).

414 *published at <http://chao.stat.nthu.edu.tw>. [accessed 27 March 2013]*

415

416 Crawford, N.G. (2010) SMOGD: software for the measurement of genetic diversity.

417 *Molecular Ecology Resources*, **10**, 556–557.

418

419 Dragulescu A.A. (2012) *xlsx: Read, write, format Excel 2007 and Excel 97/2000/XP/2003 files.*

420 <http://cran.r-project.org/web/packages/xlsx/> [accessed 27 March 2013]

421

422 du Prel, J.-B., Hommel, G., Röhrig B. & Blettner M. (2009) Confidence interval or p-value?.

423 *Deutsches Ärzteblatt international*, **106**, 335–339.

424

425 Excoffier, L. & Heckel, G. (2006) Computer programs for population genetics data analysis: a  
426 survival guide. *Nature Reviews. Genetics*, **7**, 745–58.

427

428 Gaile D.P., Shepherd L.A., Sucheston L., Bruno, A. & Manly K.F.(2012) *sendplot: Tool for*

429 *sending interactive plots with tool-tip content.* <http://cran.r->

430 [project.org/web/packages/sendplot/](http://cran.r-project.org/web/packages/sendplot/) [accessed 27 March 2013]

431

432 Gerlach, G., Jueterbock, A., Kraemer, P., Deppermann, J. & Harmand, P. (2010) Calculations

433 of population differentiation based on  $G_{ST}$  and  $D$ : forget  $G_{ST}$  but not all of statistics!.

434 *Molecular Ecology*, **19**, 3845–3852.

435

436 Goudet, J. (2004) hierfstat, a package for R to compute and test hierarchical F-statistics.  
437 *Molecular Ecology Notes*, **5**, 184–186.

438

439 Hedrick, P.W. (2005) A standardized genetic differentiation measure. *Evolution*, **59**, 1633–  
440 1638.

441

442 Hedrick, P.W. (1999) Highly variable loci and their interpretation in evolution and  
443 conservation. *Evolution*, **53**, 313–318.

444

445 Jost, L. (2008)  $G_{ST}$  and its relatives do not measure differentiation. *Molecular Ecology*, **17**,  
446 4015–4026.

447

448 Karl S.A., Toonen R.J., Grant W.S. & Bowen B.W. (2012) Common misconceptions in  
449 molecular ecology: echoes of the modern synthesis. *Molecular Ecology*, **21**, 4171–4189.

450

451 Keenan, K., Bradley, C.R., Magee, J.J., Hynes, R.A., Kennedy, R.J., Crozier, W.W., Poole, R.,  
452 Cross, T.F., McGinnity, P. & Prodöhl, P.A. (in press) Beaufort Trout MicroPlex: A high  
453 throughput multiplex platform comprising 38 informative microsatellite loci for use in  
454 brown trout and sea trout (*Salmo trutta* L.) genetics studies. *Journal of Fish Biology* (in  
455 press).

456

457 Lemon, J. (2006) Plotrix: a package in the red light district of R. *R News*, **6**, 8–12.

458

459 Meirmans, P.G. & Van Tienderen, P.H. (2004) GENOTYPE and GENODIVE: two programs for

460 the analysis of genetic diversity of asexual organisms. *Molecular Ecology Notes*, **4**, 792–794.

461

462 Meirmans, P.G. & Hedrick, P.W. (2011) Assessing population structure:  $F_{ST}$  and related  
463 measures. *Molecular Ecology Resources*, **11**, 5-18.

464

465 Nei, M. & Chesser, R.K. (1983) Estimation of fixation indices and gene diversities. *Annals of*  
466 *Human Genetics*, **47**, 253–259.

467

468 R Development Core Team (2012) R: A Language and Environment for Statistical Computing.

469 R Foundation for Statistical Computing, Vienna. <http://www.R-project.org> [accessed 27  
470 March 2013]

471

472 Raymond, M. & Rousset, F. (1995) GENEPOP (version 1.2): population genetics software for  
473 exact tests and ecumenicism. *Journal of Heredity*, **86**, 248.

474

475 Revolution Analytics (2012a) *doParallel: Foreach parallel adaptor for the parallel package*.

476 <http://cran.r-project.org/web/packages/doParallel/> [accessed 27 March 2013]

477

478 Revolution Analytics (2012b) *foreach: Foreach looping construct for R*. <http://cran.r->

479 [project.org/web/packages/foreach/](http://cran.r-project.org/web/packages/foreach/) [accessed 27 March 2013]

480

481 Revolution Analytics (2012c) *iterators: Iterator construct for R*. <http://cran.r->

482 [project.org/web/packages/iterators/](http://cran.r-project.org/web/packages/iterators/) [accessed 27 March 2013]

483

484 Rosenberg, N.A., Li, L.M., Ward, R. & Pritchard, J.K. (2003) Informativeness of genetic  
485 markers for inference of ancestry. *The American Journal of Human Genetics*, **73**, 1402–1422.  
486

487 RStudio and Inc. (2012) *shiny: Web Application Framework for R*. [http://cran.r-](http://cran.r-project.org/web/packages/shiny/)  
488 [project.org/web/packages/shiny/](http://cran.r-project.org/web/packages/shiny/) [accessed 27 March 2013]  
489

490 Skrbinišek, T., Jelenčič, M., Waits, L.P., Potočnik, H., Kos, I. & Trontelj, P. (2012) Using a  
491 reference population yardstick to calibrate and compare genetic diversity reported in  
492 different studies: an example from the brown bear. *Heredity*, **109**, 299–305.  
493

494 Wagenmakers, E.J. (2007) A practical solution to the pervasive problems of  $p$ -values.  
495 *Psychonomic Bulletin & Review*, **14**, 779–804.  
496

497 Weir, B.S. & Cockerham, C.C. (1984) Estimating F-statistics for the analysis of population  
498 structure. *Evolution*, **38**, 1358–1370.  
499

500 Whitlock, M.C. (2011)  $G'_{ST}$  and  $D$  do not replace  $F_{ST}$ . *Molecular Ecology*, **20**, 1083–91.  
501

502 Winter, D.J. (2012) mmod: an R library for the calculation of population differentiation  
503 statistics. *Molecular Ecology Resources*, **12**, 1158-1160.  
504

505 **Table 1:** Functions of the diveRsity package  
 506

Function	Returned objects	Description
<code>chiCalc</code>	R character matrix, optional <code>.txt</code> file	Test for genetic heterogeneity between population samples using the chi-square distribution. The function provides the unique option to disregard alleles of very low frequencies using the argument <code>minFreq</code> .
<code>corPlot</code>	R graphics plot (not automatically written to file)	Correlation plotting of diversity statistics against the number of alleles per locus. The function is intended to aid in the assessment of marker suitability for the estimation of geneflow.
<code>divPart</code>	<code>.html</code> , <code>.png</code> , <code>.txt</code> , <code>.xlsx</code> , R data object	A function for the calculation of diversity partition statistics and their associated variance through bootstrapping. Global, locus and pairwise levels are addressed.
<code>divOnline</code>	NA	This function launches the web app version of <code>divPart</code> . Local resources are used when running analyses. The system default web browser is used to host the application
<code>difPlot</code>	<code>.html</code> , <code>.png</code>	Provides visualization and exploration of pairwise genetic differentiation. The function is particularly useful for data sets containing a large number of population samples.
<code>inCalc</code>	<code>.png</code> , <code>.txt</code> , <code>.xlsx</code> , R data object	A function for the calculation of allele and locus informativeness for the inference of ancestry. Bootstrap confidence intervals are also calculated.
<code>readGenepop</code>	R data object	A general purpose function designed to calculate basic descriptive parameters from raw genetic data. This function is intended as a tool for developers of population genetics software in R.
<code>divRatio</code>	R data object, <code>.txt</code> , or <code>.xlsx</code>	This function calculates the diversity ratio statistics presented in (Skrbinšek <i>et al.</i> , 2012).
<code>bigDivPart</code>	R data object, <code>.txt</code> , or <code>.xlsx</code>	This function is identical to <code>divPart</code> except for its lack of bootstrapping functionality. It is coded in a specific way to allow the sequential analysis of large number of markers (e.g. <100,000).
<code>fstOnly</code>	R data object, <code>.txt</code> , or <code>.xlsx</code>	This function calculates only Weir & Cockerham's 1984 F-statistics. The function is slightly faster than <code>divPart</code> which also calculates these statistics.
<code>divBasic</code>	R data object, <code>.txt</code> , or <code>.xlsx</code>	This function calculates basic population bases statistics such as Allelic richness, Hardy-Weinberg equilibrium and locus expected and observed heterozygosities.

508 **Table 2:** Additional packages used by the diveRcity package, along with their  
509 implementations.

Package	Implementation	Status	Citation
Xlsx	Used in <code>divPart</code> and <code>inCalc</code> to return multi-sheet <code>.xlsx</code> workbooks	Suggested	(Dragulescu, 2012)
sendplot	Used in <code>divPart</code> , <code>divPlot</code> and <code>inCalc</code> to produce tool tips for data visualisation	Suggested	(Gaile <i>et al.</i> , 2012)
doParallel	Used in <code>divPart</code> and <code>inCalc</code> for parallel computation	Suggested	(Revolution Analytics, 2012a)
parallel	Used in <code>divPart</code> and <code>inCalc</code> for parallel computation	Suggested	(R Development Core Team, 2012)
foreach	Used in <code>divPart</code> and <code>inCalc</code> for parallel computation	Suggested	(Revolution Analytics, 2012b)
iterators	Used in <code>divPart</code> and <code>inCalc</code> for parallel computation	Suggested	(Revolution Analytics, 2012c)
plotrix	Used in <code>divPlot</code> for additional plotting features	Dependency	(Lemon, 2006)
shiny	Used to build and run the web app version of the <code>divPart</code> function	Dependency	(RStudio & Inc., 2012)

510

511



512 **Figures**

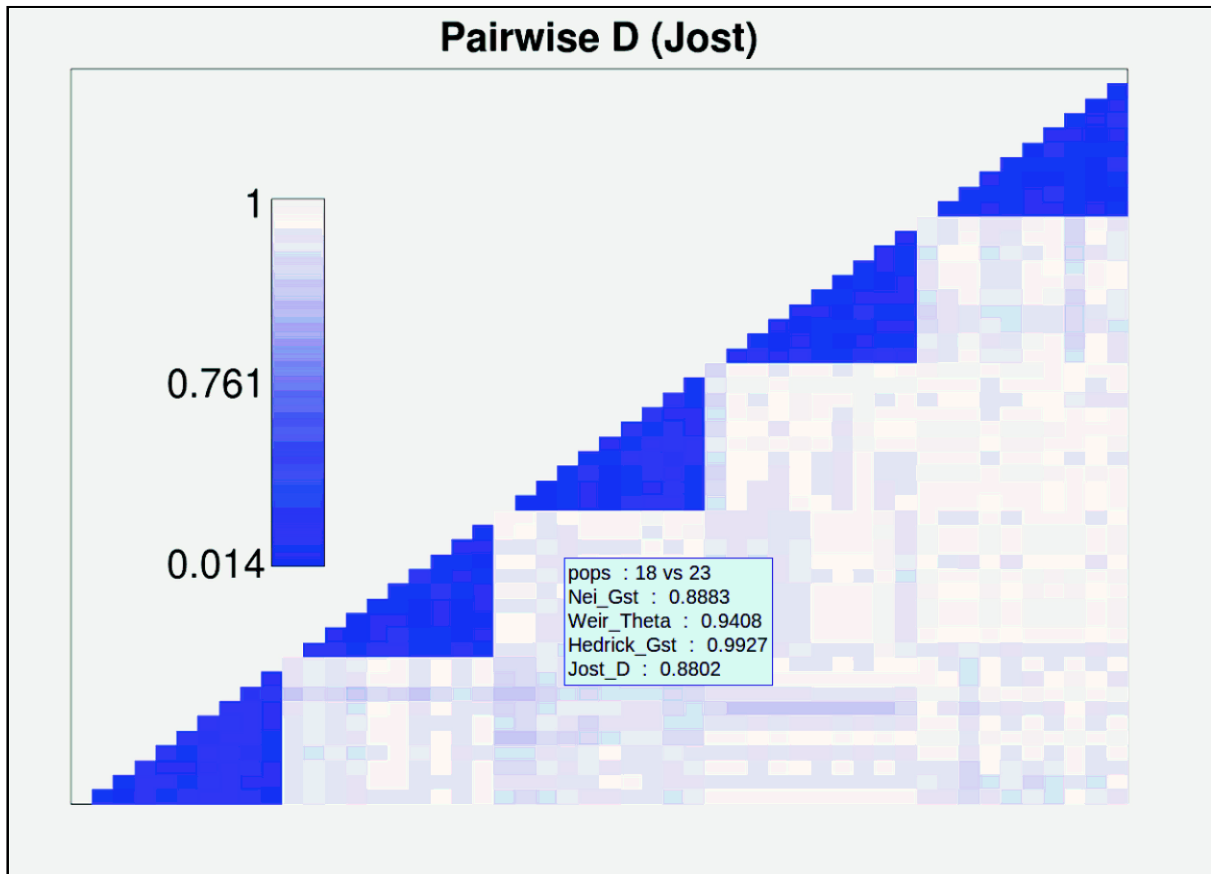
pop1	Locus1	Locus2	Locus3	Locus4	Locus5	Overall
N	46	47	47	47	45	46.4
A	4	3	11	6	19	43
%	80	100	78.57	75	73.08	81.33
<u>Ar</u>	3.57	2.92	10.36	5.41	17.12	7.88
Ho	0.67	0.57	0.87	0.64	0.76	0.7
He	0.66	0.53	0.83	0.67	0.92	0.72
HWE	0.6321	0.794	0.7286	0.8701	0.012	0.1064
pop2	Locus1	Locus2	Locus3	Locus4	Locus5	Overall
N	40	42	42	42	42	41.6
A	5	2	13	7	20	47
%	100	66.67	92.86	87.5	76.92	84.79
<u>Ar</u>	4.86	2	10.84	5.89	17.27	8.17
Ho	0.65	0.48	0.74	0.52	0.9	0.66
He	0.66	0.5	0.79	0.71	0.92	0.72
HWE	0.534	0.7642	0.9999	0.5327	0.9249	0.9988
pop3	Locus1	Locus2	Locus3	Locus4	Locus5	Overall
N	41	41	41	40	39	40.4
A	5	2	10	4	14	35
%	100	66.67	71.43	50	53.85	68.39
<u>Ar</u>	4.62	2	8.82	4	12.41	6.37
Ho	0.73	0.39	0.71	0.7	0.87	0.68
He	0.71	0.5	0.8	0.7	0.87	0.71
HWE	0	0.1604	0.9874	0.9841	0.9929	0.8389

513

514 **Figure 1.** A screen-shot of the results output format from the function `divBasic`. This table  
515 format is commonly seen in journal articles when presenting basic population genetic  
516 parameters. However, the parameters often have to be calculated in separate software  
517 packages and tabulated by authors. `diversity` aims to reduce this requirement for authors.  
518 The parameter calculated in this table are;  $N$  = Number of individuals per population sample  
519 genotyped per locus,  $A$  = Total number of alleles observed per population sample per locus,  
520 % = Percentage of total alleles observed across population samples per population sample

521 per locus,  $A_r$  = Allelic richness per locus,  $H_o$  = observed heterozygosity per locus,  $H_e$  =  
522 expected heterozygosity per locus,  $HWE$  = Hardy-Weinberg Equilibrium  $p$ -value from the  $\chi^2$   
523 goodness-of-fit tests per locus.

524

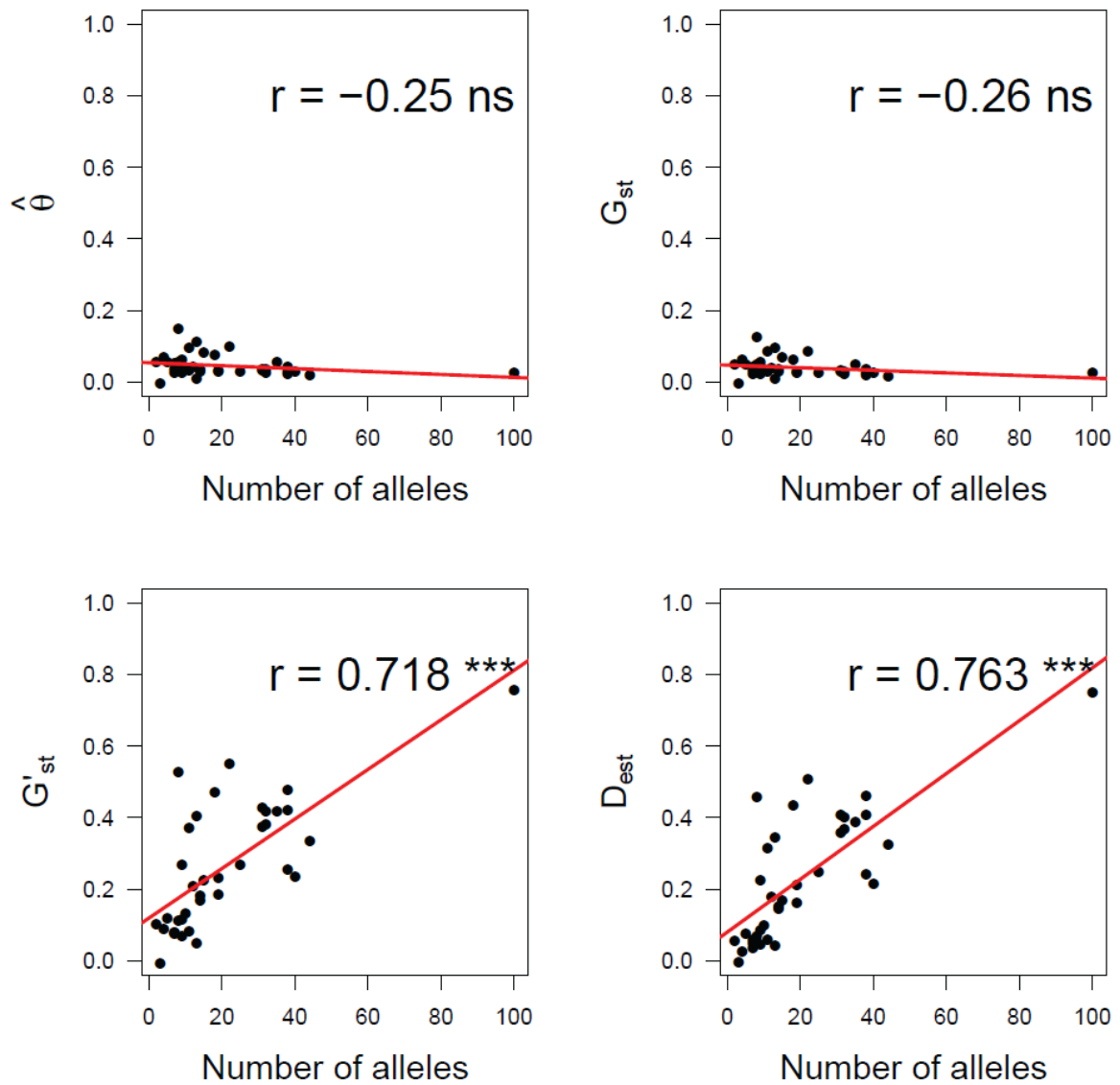


525

526

527 **Figure 2.** Visualisation of pairwise  $D_{Jost}$  (estimator), for  $N = 50$  populations. Total  
528 pairwise comparisons = 1225. This figure is returned from the `difPlot` function, which will  
529 plot diversity partitioning and differentiation estimators returned by `divPart`. Regions of  
530 dark blue represent low genetic differentiation, while light blue/white represents high  
531 differentiation. The text box caption is an example of the tool-tip information associated with  
532 each pairwise population comparison.

533



534

535

536 **Figure 3.** Correlation assessment of locus estimators  $\theta$ ,  $G_{ST}$ ,  $G'_{ST}$ , and  $D_{est}$   
 537 ( $D_{Jost}$  unbiased estimator), with locus polymorphism (total number of alleles), returned from  
 538 the `corPlot` function. Red lines represent the line of best fit and r values are Pearson product  
 539 moment correlation coefficients.