

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23

MR. BRENDAN F WRINGE (Orcid ID : 0000-0002-9482-5534)

DR. ERIC C ANDERSON (Orcid ID : 0000-0003-1326-0840)

Article type : Resource Article

hybriddetective: a workflow and package to facilitate the detection of hybridization using genomic data in R

Brendan F. Wringe^{1*}, Ryan R. E. Stanley¹, Nicholas W. Jeffery¹, Eric C. Anderson², and
Ian R. Bradbury¹

¹ Science Branch, Department of Fisheries and Oceans Canada, 80 East White Hills
Road, St. John's NL, A1C 5X1

² Fisheries Ecology Division, National Oceanic and Atmospheric Administration
Southwest Fisheries Science Center, Santa Cruz, CA, 95060

*Corresponding author, bwringe@gmail.com

Running title: hybriddetective: hybrid detection workflow

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: 10.1111/1755-0998.12704](https://doi.org/10.1111/1755-0998.12704)

This article is protected by copyright. All rights reserved

24 **Keywords:** hybrid, introgression, population genetics, population structure,
25 assignment tests, simulation

26

27 **Abstract**

28 The ability to detect and characterize hybridization in nature has long been of
29 interest to many fields of biology and often has direct implications for wildlife
30 management and conservation. The capacity to identify the presence of
31 hybridization, and quantify the numbers of individuals belonging to different hybrid
32 classes, permits inference on the magnitude of, and time scale over which,
33 hybridization has been, or is occurring. Here we present an R package and
34 associated workflow developed for the detection, with estimates of efficiency and
35 accuracy, of multi-generational hybrid individuals using genetic or genomic data in
36 conjunction with the program NEWHYBRIDS. This package includes functions for
37 the identification and testing of diagnostic panels of markers, the simulation of
38 multi-generational hybrids, and the quantification and visualization of the efficiency
39 and accuracy with which hybrids can be detected. Overall, this package delivers a
40 streamlined hybrid analysis platform, providing improvements in speed, ease of use
41 and repeatability over current *ad hoc* approaches. The latest version of the package
42 and associated documentation are available on GitHub
43 (<https://github.com/bwringe/hybriddetective>).

44 **Introduction**

45

46 Detecting and elucidating patterns of hybridization between individuals from
47 genetically distinct populations is of interest in many fields of biology (Abbott *et al.*
48 2013; Payseur & Rieseberg 2016; Todesco *et al.* 2016). Naturally occurring hybrid
49 zones - areas where genetically distinct populations come into contact and create
50 genetically (ad)mixed offspring - are important natural laboratories to study of the

51 interplay between selection and recombination (Barton & Hewitt 1985; Burke &
52 Arnold 2001; Hilbish *et al.* 2012). These areas have provided opportunities to glean
53 information to further model, and test hypotheses related to speciation (Abbott *et al.*
54 2013; Barton 2013; Dowling & Secor 1997) and the maintenance of reproductive
55 barriers (Albrechtová *et al.* 2012; Griebel *et al.* 2015; Landry *et al.* 2007), natural
56 selection (Johnson *et al.* 2010; Pruvost *et al.* 2013), and genetic recombination.
57 Hybridization can also have conservation, regulatory, and legal ramifications related
58 to the genetic structure and integrity of populations (Allendorf *et al.* 2004; Benson
59 *et al.* 2014; Boyer *et al.* 2008; Fitzpatrick *et al.* 2015; Rostgaard Nielsen *et al.* 2016),
60 or the introgression of domesticated (Fraser *et al.* 2010; Kidd *et al.* 2009; Noren *et*
61 *al.* 2005) or transgenic (Oke *et al.* 2013; Warwick *et al.* 2003) alleles into wild
62 populations.

63 In some cases, hybrid individuals can be identified morphologically (de
64 Oliveira *et al.* 2002; Ross & Cavender 1981; Solomon & Child 1978), however
65 morphological classification is notoriously imperfect (Baumsteiger *et al.* 2005;
66 Esquer-Garrigos *et al.* 2015; Hardig *et al.* 2000; Neff & Smith 1979) and does not
67 allow for the classification of hybrid category (Lamb & Avise 1987) or the
68 examination of the effect of genetic dosage (Kierzkowski *et al.* 2011; Rieseberg
69 1995). In contrast, the use of Mendelian genetic markers affords researchers the
70 ability to not only identify individuals as hybrid or purebred, but also to characterize
71 them to specific hybrid classes (e.g. pure, F₁, F₂ and backcrosses). This ability to
72 quantify the types, and numbers of individuals of different hybrid classes present,
73 allows inferences to be made on the magnitude of, and time scale over which,
74 hybridization has been, or is occurring (Anderson & Thompson 2002; Brown *et al.*
75 2004; Godinho *et al.* 2015; Saarman & Pogson 2015).

76 Many statistical approaches have been put forward to use genetic markers to
77 identify hybrids (Anderson 2009), and some of these have been incorporated into
78 widely used, and cited software programs (e.g. NEWHYBRIDS [Anderson &
79 Thompson 2002]; STRUCTURE [Hubisz *et al.* 2009]; GENODIVE [Meirmans & Van
80 Tienderen 2004]). However, the analyses conducted by these programs is but one
81 step in the path to go from individual genotypes, to the detection and assignment of

82 those individuals to a hybrid class, with quantifiable levels of certainty. The process
83 of performing hybrid analyses currently entails the use of multiple, standalone
84 programs, many of which require data to be provided in a unique format (Lischer &
85 Excoffier 2012; Stanley *et al.* 2017). Furthermore, the reliance on the user for file
86 management, and for manually implementing individual analyses with separate
87 programs in addition to affording opportunity for human error, leads to a disjunct
88 analytical process with a steep learning curve that lacks the efficiency and
89 repeatability of a true workflow.

90 Here we describe the R package *hybriddetective* and associated workflow for
91 hybrid identification developed in the R computer language (R Development Core
92 Team 2016). The package and workflow encompass every aspect of the hybrid
93 identification procedure. Specifically, we include functions for (1) panel design, and
94 the quantification of the efficiency, accuracy and power of panels of diagnostic
95 markers; (2) error checking and diagnostics; and (3) quantification, and
96 visualization of accuracy and assignment power of the selected panel(s).
97 *hybriddetective's* simulation and panel selection functions have been designed to
98 work in concert, as a workflow, to improve the accuracy, and reduce the
99 overestimation of assignment certainty (Anderson & Thompson 2002), and
100 concomitantly reduce high-grading bias (described in detail below; Anderson 2010).
101 This package alleviates much of the complexity in the hybrid detection process,
102 reduces the potential for human error, and at the same time offers significant speed
103 improvements over previous *ad hoc* methodologies.

104 **Description of the package**

105
106 *hybriddetective* is compiled as an R package which facilitates a workflow within the
107 R environment for the detection of hybrids based on genotypic/genomic
108 information using the program NEWHYRIDS (Anderson & Thompson 2002), and
109 provides a comprehensive and repeatable framework to move from genotypic data
110 to the identification, with quantifiable certainty, of hybrid individuals.

111 *hybriddetective* is comprised of 14 functions (Table 1), three example datasets, and
112 a README. Function descriptions (Table 1), example data, and installation
113 instructions are available online <https://github.com/bwringe/hybriddetective>. For
114 an example of the *hybriddetective* workflow, see Jeffery *et al.* (2017), and
115 Supplementary Figure 1 . We chose to implement hybrid detection using the
116 program NEWHYBRIDS (Anderson & Thompson 2002) because it permits the
117 assignment of individuals to hybrid class (i.e. pure-bred, F₁, F₂, and back-crosses)
118 and does not require *a priori* knowledge of the allele frequencies of the two
119 populations being tested (Anderson & Thompson 2002). Moreover, NEWHYBRIDS is
120 widely used, having been cited over 800 times as of the time of this writing.

121 **Description of the workflow**

122
123 The workflow can be broken down into three major elements: 1) data preparation,
124 2) error checking and diagnostics, and 3) quantification and analysis. **Data**
125 **preparation** encompasses the process of selecting the *n* most informative loci from
126 amongst the genotypic data available, and the simulation of multi-generational
127 hybrids. After analyzing the simulated data with NEWHYBRIDS, **error checking and**
128 **diagnostics** functions confirm that NEWHYBRIDS MCMC chains reached
129 convergence and **quantification and analysis** functions test, quantify, and visualize
130 the accuracy and assignment power of the selected panel(s). The workflow and the
131 functions used in each step are illustrated and described in in Figure 1, and Table 1,
132 respectively. We have also included a brief section on the implementation of
133 (parallel) NEWHYBRIDS analyses using the related R package *parallelenewhybrid*
134 (Wringe *et al.* 2017).

135 **Data preparation**

136 *Panel selection*

137

138 Panel selection is the process of selecting from amongst the available markers (i.e.
139 thousands to several hundreds of thousands as produced by RAD-seq) a subset that
140 together permit accurate identification of hybrids. In our workflow, the function
141 *getTopLoc* is used to develop a panel of user defined size, of the most informative
142 (based on global Weir and Cockerham (1984)'s F_{ST}) loci that are not in linkage
143 disequilibrium (LD). Genotype data of individuals known (or suspected with high
144 certainty (Oliveira *et al.* 2015)), to be of pure ancestry from the two populations
145 potentially hybridizing are used as input for *getTopLoc*. *getTopLoc* first randomly
146 creates two subsamples, each comprised of 50% of the individuals from each of the
147 two populations, to create validation and training datasets. To prevent any "high-
148 grading" bias (i.e. upward bias in the estimation of predictive capacity caused when
149 the same data is used to both select and validate panels of markers), *getTopLoc* uses
150 subsampling to ensure the same individuals are not used to create the panel and to
151 validate it. The function uses the training dataset to calculate the global, locus-
152 specific Weir and Cockerham's F_{ST} and ranks loci by this metric. Pairwise LD is then
153 calculated using the training dataset for all loci within one or both populations at the
154 users' discretion. During this process the r^2 threshold above which to consider a pair
155 of loci to be in LD can be defined by the user. Any loci that are in LD are removed,
156 because NEWHYBRIDS assumes no linkage, and each locus is treated as
157 independent. *getTopLoc* returns a list of panel loci names, a list of individuals (IDs)
158 in the validation dataset, and the genotypes of those individuals at the panel loci.

159 Importantly, random sampling selects the individuals in the training and
160 validation datasets, so the individuals and corresponding panel can vary each time
161 the function is run. The variance in global pairwise F_{ST} , and hence the loci returned
162 between runs, will likely be greatest where sample sizes for one or both populations
163 are small, and consequently subsampling is more apt to impart stochastic variances
164 in allele/gene frequencies.

165 *Construction of multi-generational simulated hybrids*

166

167 The next step in our workflow is to generate simulated multi-generational hybrid
168 datasets using the genotypic data from the validation dataset exported by
169 *getTopLoc*. The two simulation functions, *freqbasedsim_GTFreq* and
170 *freqbasedsim_AlleleSample* differ in the way in which they create hybrids.
171 *freqbasedsim_GTFreq* was designed to simulate individuals within the R
172 environment analogously to the commonly used hybrid simulation program
173 HYBRIDLAB (Nielsen *et al.* 2006). In *freqbasedsim_GTFreq*, like in HYBRIDLAB,
174 individuals in generation $t+1$ are created by sampling one allele per locus from the
175 generation t parental populations, based on the allele frequencies in either
176 population. Unlike HYBRIDLAB, *freqbasedsim_GTFreq* creates multi-generational
177 hybrids, each time it is run, and requires only a single data file to do so. In controlled
178 comparisons with HYBRIDLAB we find *freqbasedsim_GTFreq* to be more than 20X
179 faster when creating multiple independent simulations (See Supplemental Table 1).

180 The other hybrid simulation function, *freqbasedsim_AlleleSample*, was
181 designed with the intent of providing an additional simulation method. It first
182 randomly subsamples a proportion of individuals from each of the two populations
183 provided to it and only the alleles of these individuals will be available during the
184 subsequent simulation. Secondly, to conduct the actual simulations, each locus in
185 individuals in generation $t+1$ is simulated by randomly sampling without
186 replacement, one allele from among all the alleles present at that locus from one of
187 the parental populations at time t , then combining it with an allele chosen in the
188 same manner from the other parental population at time t . In this case, the number
189 of individuals that can be simulated in a given hybrid generation is therefore
190 dependent upon the number of individuals sampled in the first step.

191 (*Parallel*) *NEWHYBRIDS* analyses

192
193 For actual hybrid identification, we encourage users to take advantage of the R
194 package **parallelnewhybrid**, which was developed to run *NEWHYBRIDS* in parallel
195 thus providing significant speed improvements (Wringe *et al.* 2017). Furthermore,
196 the error checking and analytical functions described below were designed to work

197 with the file structure created by **parallelnewhybrid**. **parallelnewhybrid**, and
198 documentation describing its installation and operation can be found at
199 <https://github.com/bwringe/parallelnewhybrid>.

200

201 **Error checking and diagnostics**

202 *Check Markov chain convergence*

203

204 As with any MCMC process using Gibbs sampling, chain convergence in
205 NEWHYBRIDS is dependent upon the ‘topography’ of the probability space relative
206 to the starting point of the chain. Occasionally, the MCMC chains in NEWHYBRIDS
207 analyses will fail to converge. In these cases NEWHYBRIDS will almost invariably
208 report that (nearly) all individuals have the highest posterior probability of
209 membership in the F₂ hybrid class, a result that is clearly erroneous. To this end, the
210 function *nh_preCheckR* quickly checks the results of NEWHYBRIDS flagging those
211 that may have failed to converge, and the function *nh_multiplotR* complements it by
212 visualizing its results. *nh_preCheckR* inspects the NEWHYBRIDS output and
213 identifies the individuals that are known to be pure-bred in origin, and checks that a
214 user defined proportion of these individuals have not been assigned posterior
215 probability of assignments (PofZ; Anderson 2003) to the F₂ hybrid class in excess of
216 a user defined threshold. If these conditions are violated, the user is prompted to
217 verify the(se) result(s). *nh_multiplotR* permits the user to visualize the cumulative
218 posterior probability of assignment for all genotype frequency classes for each
219 individual. *nh_multiplotR* can thus be used to confirm and compliment the results of
220 *nh_preCheckR*, as well as quickly visualize the results of multiple NEWHYBRIDS
221 analyses.

222 **Quantification and analysis**

223 *Assess panel accuracy*

224

225 The next step in the workflow, after confirming convergence, is to assess the ability
226 of NEWHYBRIDS to assign simulated individuals of known hybrid ancestry to the
227 correct genotype frequency class given the genotypes of the individuals at the loci in
228 the selected panel. Because it is impossible to statistically validate the assumed
229 distribution of priors, and the efficacy of the loci in a panel *a priori* (Anderson 2003;
230 Nielsen *et al.* 2006; Oliveira *et al.* 2008), simulations are often employed to evaluate
231 power (Anderson 2003; Nielsen *et al.* 2006; Vähä & Primmer 2006). Also, Anderson
232 and Thompson (2002), note that the power of NEWHYBRIDS to distinguish among
233 genotype frequencies classes will vary across classes. Thus when evaluating a
234 potential threshold value of posterior probability of assignment for assigning
235 genotype frequency class membership, the effect of choice of posterior probability
236 of assignment value on efficiency, accuracy and overall performance (Vähä &
237 Primmer 2006), as well as on both Type I and Type II error should be considered
238 simultaneously for each genotype frequency class, and for the differentiation of
239 purebreds from any type of hybrid.

240 In order to allow researchers to better evaluate the effect of choice of critical
241 posterior probability of assignment threshold (i.e. posterior probability value above
242 which assignment to a given hybrid class is accepted) on assignment success, we
243 have developed the function *hybridPowerComp*. *hybridPowerComp* calculates the
244 number of individuals of known hybrid class correctly assigned over the total
245 number of individuals known to belong to that class for posterior probability of
246 assignment thresholds between 0.50 and 1.0 (i.e. number detected / number
247 expected; "efficiency" sensu Vähä & Primmer 2006). This is done for each hybrid
248 frequency class (Figure 2), as well as separately for the two parental classes, and all
249 hybrids classes considered together (i.e. posterior probability of assignment for
250 hybrid is the sum of all of F₁, F₂, BC1, BC2). In addition, *hybridPowerComp*
251 calculates and plots the number of individuals correctly assigned to a class over the
252 total number of individuals assigned to that class (i.e. "accuracy" sensu Vähä &
253 Primmer 2006)(Figure 3), and the "power" (i.e. the product of "efficiency" and
254 "accuracy" sensu Vähä & Primmer 2006) of the panel . Similarly, the number of
255 individuals wrongly deemed to belong to hybrid genotype frequency classes divided

256 by the total number of known pure individuals (i.e. type I error; Burgarella *et al.*
257 2009), and the proportion of individuals misclassified (i.e. type II error) are
258 assessed and plotted. *hybridPowerComp* allows visualization of the distribution of
259 posterior probability of assignment values by plotting them for each genotype
260 frequency class, as well as for all hybrid classes considered together (refer to
261 Supplementary Table 2 for a list of the plots produced by *hybridPowerComp*).

262 ▪ The function *nh_panel_delta_plotR* complements *hybridPowerComp* by
263 visualizing the efficacy of different panel sizes for each genotype frequency class and
264 can be used during the assessment of panel accuracy phase of the workflow.

265

266 *Combine simulated and experimental data for analysis*

267

268 Once the panel and critical posterior probability of assignment threshold(s) have
269 been finalized, the experimental/unknown data can be analyzed. Combining
270 simulated data with the unknown/experimental data (1) assists with the
271 interpretation of results in the absence of known individuals, and (2) allows the user
272 the option to designate the genotype frequency class membership of known
273 individuals, to improve assignment power (Anderson 2003; Anderson & Thompson
274 2002).

275 The function *nh_analysis_generateR* allows researchers to specify both the
276 unknown and experimental genotype data to analyze and the simulated data to
277 combine with it, thus facilitating reproducibility of analyses as well as the ability to
278 use the same simulated dataset(s) from which the critical posterior probability of
279 assignment values were determined. The function *nh_analysis_simulateR_generateR*
280 permits users to quickly create analysis-ready datasets when panel development
281 and/or more conservative simulation methodology are not required. This function
282 uses the frequency based simulation algorithm and simulation options of
283 *freqbasedsim_GTFreq* to create simulated hybrids based on supplied genotype data,
284 and then merge them with experimental or unknown genotypes.

285 **Conclusions**

286

287 Here we have shown that the use of *hybriddetective* as part of a workflow in the
288 detection of hybrids has clear and quantifiable benefits over the generally *ad hoc*,
289 methods normally used. *hybriddetective* provides researchers an efficient platform
290 for reproducible analyses of hybridization within the R computational language.
291 Furthermore, the interoperability of *hybriddetective* for the simulation of multi-
292 generation hybrid datasets and the separate R package *parallelnewhybrid* (Wringe
293 *et al.* 2017) to efficiently and automatically execute runs of NEWHYBRIDS in
294 parallel, makes it tractable to quantify the expected variability in hybrid assignment
295 success.

296 In conclusion, we have created an R package and associated workflow for the
297 detection, with quantifiable accuracy, efficiency and power, of multi-generational
298 hybrid individuals using genetic or genomic data with the program NEWHYBRIDS.
299 This package includes functions for the development and testing of diagnostic
300 panels of markers, the simulation of multi-generational hybrids, and the
301 quantification and visualization of the accuracy with which (simulated) hybrids can
302 be detected. Use of this package offers improvements in the repeatability, speed, and
303 ease of use over conventional approaches.

304 **Acknowledgements**

305

306 The authors wish to thank Marion Sinclair-Waters, Justine [Létourneau](#), and Anne-
307 [Laure Ferchaud](#) for their help bug-checking the code. We also thank Thierry
308 Gosselin for encouraging us to publish this package. The manuscript was greatly
309 improved by comments from Sarah Lehnert and three anonymous reviewers. This
310 work was supported by a Natural Sciences and Engineering Research Council
311 Strategic Project Grant, a Natural Sciences and Engineering Research Discovery
312 Grant, and Canadian Healthy Oceans Network, and Fisheries and Oceans Canada

313 funding (International Governance Strategy; Programme for Aquaculture
314 Regulatory Research; Genomics Research and Development Initiative) to I.R.B.

315 **Author contributions**

316 B.F.W. wrote the manuscript and the package code, and developed the supporting
317 documentation and example data files hosted on GitHub. R.R.E.S, N.F.W., E.C.A., and
318 I.R.B. all contributed to the initial concept, development of the code, and associated
319 documentation, as well as assisting in writing of the manuscript.
320

321 **Data Accessibility**

322 The package, user manual, README, example workflow, and example data sets are
323 all available online from <https://github.com/bwringe/hybriddetective>.

324 **References**

325

326 Abbott R, Albach D, Ansell S, *et al.* (2013) Hybridization and speciation. *Journal of*
327 *Evolutionary Biology* **26**, 229-246.

328 Albrechtová J, Albrecht T, Baird SJE, *et al.* (2012) Sperm-related phenotypes
329 implicated in both maintenance and breakdown of a natural species barrier
330 in the house mouse. *Proceedings of the Royal Society B-Biological Sciences*
331 **279**, 4803-4810.

332 Allendorf FW, Leary RF, Hitt NP, *et al.* (2004) Intercrosses and the US Endangered
333 Species Act: Should hybridized populations be included as westslope
334 cutthroat trout? *Conservation Biology* **18**, 1203-1213.

335 Anderson EC (2003) User's guide to the program NewHybrids Version 1.1 beta.
336 Department of Integrative Biology, University of California, Berkeley,
337 Berkeley, California.

338 Anderson EC (2009) Statistical methods for identifying hybrids and groups. In:
339 *Population genetics for animal conservation* (eds. Bertorelle G, Bruford MW,
340 Hauff HC, Rizzoli A, Vernesi C), pp. 25-41. Cambridge University Press, New
341 York.

342 Anderson EC (2010) Assessing the power of informative subsets of loci for
343 population assignment: standard methods are upwardly biased. *Molecular*
344 *Ecology Resources* **10**, 701-710.

345 Anderson EC, Thompson EA (2002) A model-based method for identifying species
346 hybrids using multilocus genetic data. *Genetics* **160**, 1217-1229.

347 Barton NH (2013) Does hybridization influence speciation? *Journal of Evolutionary*
348 *Biology* **26**, 267-269.

349 Barton NH, Hewitt GM (1985) Analysis of hybrid zones. *Annual Review of Ecology*
350 *and Systematics* **16**, 113-148.

351 Baumsteiger J, Hankin D, Loudenslager EJ (2005) Genetic analyses of juvenile
352 steelhead, coastal cutthroat trout, and their hybrids differ substantially from

353 field identifications. *Transactions of the American Fisheries Society* **134**, 829-
354 840.

355 Benson JF, Patterson BR, Mahoney PJ (2014) A protected area influences genotype-
356 specific survival and the structure of a *Canis* hybrid zone. *Ecology* **95**, 254-
357 264.

358 Boyer MC, Muhlfeld CC, Allendorf FW (2008) Rainbow trout (*Oncorhynchus mykiss*)
359 invasion and the spread of hybridization with native westslope cutthroat
360 trout (*Oncorhynchus clarkii lewisi*). *Canadian Journal of Fisheries and*
361 *Aquatic Sciences* **65**, 658-669.

362 Brown KH, Patton SJ, Martin KE, *et al.* (2004) Genetic analysis of interior Pacific
363 Northwest *Oncorhynchus mykiss* reveals apparent ancient hybridization
364 with westslope cutthroat trout. *Transactions of the American Fisheries*
365 *Society* **133**, 1078-1088.

366 Burgarella C, Lorenzo Z, Jabbour-Zahab R, *et al.* (2009) Detection of hybrids in
367 nature: application to oaks (*Quercus suber* and *Q. ilex*). *Heredity* **102**, 442-
368 452.

369 Burke JM, Arnold ML (2001) Genetics and the fitness of hybrids. *Annual Review of*
370 *Genetics* **35**, 31-52.

371 de Oliveira AC, Garcia AN, Cristofani M, Machado MA (2002) Identification of citrus
372 hybrids through the combination of leaf apex morphology and SSR markers.
373 *Euphytica* **128**, 397-403.

374 Dowling TE, Secor CL (1997) The role of hybridization and introgression in the
375 diversification of animals. *Annual Review of Ecology and Systematics* **28**,
376 593-619.

377 Esquer-Garrigos Y, Hugueny B, Ibañez C, *et al.* (2015) Detecting natural
378 hybridization between two vulnerable Andean pupfishes (*Orestias agassizii*
379 and *O. luteus*) representative of the Altiplano endemic fisheries.
380 *Conservation Genetics* **16**, 717-727.

381 Fitzpatrick BM, Ryan ME, Johnson JR, Corush J, Carter ET (2015) Hybridization and
382 the species problem in conservation. *Current Zoology* **61**, 206-216.

383 Fraser DJ, Minto C, Calvert AM, Eddington JD, Hutchings JA (2010) Potential for
384 domesticated-wild interbreeding to induce maladaptive phenology across
385 multiple populations of wild Atlantic salmon (*Salmo salar*). *Canadian Journal*
386 *of Fisheries and Aquatic Sciences* **67**, 1768-1775.

387 Godinho R, López-Bao JV, Castro D, *et al.* (2015) Real-time assessment of
388 hybridization between wolves and dogs: combining noninvasive samples
389 with ancestry informative markers. *Molecular Ecology Resources* **15**, 317-
390 328.

391 Griebel J, Giessler S, Poxleitner M, *et al.* (2015) Extreme environments facilitate
392 hybrid superiority - the story of a successful *Daphnia galeata x longispina*
393 hybrid clone. *PLoS One* **10**, e0140275.

394 Hardig TM, Brunsfeld SJ, Fritz RS, Morgan M, Orians CM (2000) Morphological and
395 molecular evidence for hybridization and introgression in a willow (*Salix*)
396 hybrid zone. *Molecular Ecology* **9**, 9-24.

397 Hilbish TJ, Lima FP, Brannock PM, *et al.* (2012) Change and stasis in marine hybrid
398 zones in response to climate warming. *Journal of Biogeography* **39**, 676-687.

399 Hubisz MJ, Falush D, Stephens M, Pritchard JK (2009) Inferring weak population
400 structure with the assistance of sample group information. *Molecular Ecology*
401 *Resources* **9**, 1322-1332.

402 Jeffery NW, DiBacco C, Wringe BF, *et al.* (2017) Genomic evidence of hybridization
403 between two independent invasions of European green crab (*Carcinus*
404 *maenas*) in the Northwest Atlantic. *Heredity (Edinb)*.

405 Johnson JR, Fitzpatrick BM, Shaffer HB (2010) Retention of low-fitness genotypes
406 over six decades of admixture between native and introduced tiger
407 salamanders. *BMC Evolutionary Biology* **10**, 147.

408 Kidd AG, Bowman J, Lesbarreres D, Schulte-Hostedde AI (2009) Hybridization
409 between escaped domestic and wild American mink (*Neovison vison*).
410 *Molecular Ecology* **18**, 1175-1186.

411 Kierzkowski P, Pasko L, Rybacki M, Socha M, Ogielska M (2011) Genome dosage
412 effect and hybrid morphology - the case of the hybridogenetic water frogs of
413 the *Pelophylax esculentus* complex. *Annales Zoologici Fennici* **48**, 56-66.

414 Lamb T, Avise JC (1987) Morphological variability in genetically defined categories
415 of anuran hybrids. *Evolution* **41**, 157-165.

416 Landry CR, Hartl DL, Ranz JM (2007) Genome clashes in hybrids: insights from gene
417 expression. *Heredity* **99**, 483-493.

418 Lischer HEL, Excoffier L (2012) PGDSpider: an automated data conversion tool for
419 connecting population genetics and genomics programs. *Bioinformatics* **28**,
420 298-299.

421 Meirmans PG, Van Tienderen PH (2004) GENOTYPE and GENODIVE: two programs
422 for the analysis of genetic diversity of asexual organisms. *Molecular Ecology*
423 *Notes* **4**, 792-794.

424 Neff NA, Smith GR (1979) Multivariate-analysis of hybrid fishes. *Systematic Zoology*
425 **28**, 176-196.

426 Nielsen EE, Bach LA, Kotlicki P (2006) HYBRIDLAB (version 1.0): a program for
427 generating simulated hybrids from population samples. *Molecular Ecology*
428 *Notes* **6**, 971-973.

429 Noren K, Dalen L, Kvaloy K, Angerbjorn A (2005) Detection of farm fox and hybrid
430 genotypes among wild arctic foxes in Scandinavia. *Conservation Genetics* **6**,
431 885-894.

432 Oke KB, Westley PAH, Moreau DTR, Fleming IA (2013) Hybridization between
433 genetically modified Atlantic salmon and wild brown trout reveals novel
434 ecological interactions. *Proceedings of the Royal Society B-Biological*
435 *Sciences* **280**.

436 Oliveira R, Godinho R, Randi E, Alves PC (2008) Hybridization versus conservation:
437 are domestic cats threatening the genetic integrity of wildcats (*Felis silvestris*
438 *silvestris*) in Iberian Peninsula? *Philosophical Transactions of the Royal*
439 *Society B-Biological Sciences* **363**, 2953-2961.

440 Oliveira R, Randi E, Mattucci F, *et al.* (2015) Toward a genome-wide approach for
441 detecting hybrids: informative SNPs to detect introgression between
442 domestic cats and European wildcats (*Felis silvestris*). *Heredity* **115**, 195-
443 205.

- 444 Payseur BA, Rieseberg LH (2016) A genomic perspective on hybridization and
445 speciation. *Molecular Ecology* **25**, 2337-2360.
- 446 Pruvost NBM, Hollinger D, Reyer H-U (2013) Genotype-temperature interactions on
447 larval performance shape population structure in hybridogenetic water frogs
448 (*Pelophylax esculentus* complex). *Functional Ecology* **27**, 459-471.
- 449 R Development Core Team (2016) *R: A language and environment for statistical*
450 *computing* R Foundation for Statistical Computing, Vienna, Austria.
- 451 Rieseberg LH (1995) The role of hybridization in evolution: old wine in new skins.
452 *American Journal of Botany* **82**, 944-953.
- 453 Ross MR, Cavender TM (1981) Morphological Analyses of Four Experimental
454 Intergeneric Cyprinid Hybrid Crosses. *Copeia* **1981**, 377-387.
- 455 Rostgaard Nielsen L, Brandes U, Dahl Kjaer E, Fjellheim S (2016) Introduced Scotch
456 broom (*Cytisus scoparius*) invades the genome of native populations in
457 vulnerable heathland habitats. *Molecular Ecology* **25**, 2790-2804.
- 458 Saarman NP, Pogson GH (2015) Introgression between invasive and native blue
459 mussels (genus *Mytilus*) in the central California hybrid zone. *Molecular*
460 *Ecology* **24**, 4723-4738.
- 461 Solomon DJ, Child AR (1978) Identification of juvenile natural hybrids between
462 Atlantic salmon (*Salmo salar* L) and trout (*Salmo trutta* L). *Journal of Fish*
463 *Biology* **12**, 499-&.
- 464 Stanley RR, Jeffery NW, Wringe BF, DiBacco C, Bradbury IR (2017) genepopedit: a
465 simple and flexible tool for manipulating multilocus molecular data in R.
466 *Molecular Ecology Resources* **17**, 12-18.
- 467 Todesco M, Pascual MA, Owens GL, *et al.* (2016) Hybridization and extinction.
468 *Evolutionary Applications*.
- 469 Vähä JP, Primmer CR (2006) Efficiency of model-based Bayesian methods for
470 detecting hybrid individuals under different hybridization scenarios and with
471 different numbers of loci. *Molecular Ecology* **15**, 63-72.
- 472 Warwick SI, Simard MJ, Légère A, *et al.* (2003) Hybridization between transgenic
473 *Brassica napus* L. and its wild relatives: *Brassica rapa* L., *Raphanus*

474 *raphanistrum* L., *Sinapis arvensis* L., and *Erucastrum gallicum* (Willd.) OE
475 Schulz. *Theoretical and Applied Genetics* **107**, 528-539.
476 Weir BS, Cockerham CC (1984) Estimating F-statistics for the analysis of population-
477 structure. *Evolution* **38**, 1358-1370.
478 Wringe BF, Stanley RR, Jeffery NW, Anderson EC, Bradbury IR (2017)
479 parallelnewhybrid: an R package for the parallelization of hybrid detection
480 using newhybrids. *Molecular Ecology Resources* **17**, 91-95.
481

Author Manuscript

Table 1 – Functions included in the *hybriddetective* R package, a synopsis of their purpose, and which of the three major elements they are used in.

Function Name	Synopsis	Main Use
<i>getTopLoc</i>	Creates a panel comprised of the n (user-specified) most informative (based on highest loci-specific F_{ST} s), markers not in linkage disequilibrium. The function randomly assigns half the individuals in each of the two populations to be used to calculate loci-specific Weir and Cockerham's F_{ST} ¹ , and returns the genotypes at the n loci of the other half to be used to test the efficacy of the panel to avoid high-grading bias ² .	Data Preparation
<i>freqbasedsim_GTFreq</i>	Creates simulated multi-generational (i.e. Pure 1, Pure 2, F ₁ , F ₂ , BC1, BC2) hybrids based on the allele frequencies in the two populations provided. The user can specify the number of individuals in each of the hybrid classes to be created.	Data Preparation
<i>freqbasedsim_AlleleSample</i>	Creates simulated multi-generational hybrids by randomly sampling, without replacement, two alleles per loci from a proportion of the individual genotypes provided. The user is able to specify the proportion of genotypes to sample, as well as the number of individuals of each hybrid class to create.	Data Preparation
<i>nh_analysis_generateR</i>	Merges a file composed of simulated hybrid genotypes with a file	Data

	containing the genotypes of unknown/experimental individuals to produce a file suitable to ascertain the hybrid class of the unknowns. The user is able to specify which hybrid classes from the simulated dataset to include in the output.	Preparation
<i>nh_analysis_simulateReferenceR</i>	Creates a simulated multi-generational hybrid reference dataset from user provided data, and then merges it with the genotypes of unknown/experimental individuals. This function will create a new simulated dataset each time it is run using the same simulation methodology as <i>freqbasedsim_GTFreq</i> .	Data Preparation
<i>nh_subsetR</i>	Removes subsets of desired loci from NEWHYBRIDS formatted files so that the efficacy of panels of various sizes can be assessed.	Data Preparation
<i>nh_Zscore</i>	Allows the user to assign known hybrid category designations to individuals in NEWHYBRIDS formatted files	Data Preparation
<i>nh_preCheckR</i>	Checks all NEWHYBRIDS results within a directory and flags those that show evidence that the Markov chain may have failed to converge. This is done by evaluating the proportion of known Pure Population 1 or 2 individuals in which the posterior probability of assignment to F ₂ exceeds a threshold. The user may specify both the proportion of individuals and the PofZ threshold.	Error Checking and Diagnostics
<i>nh_multplotR</i>	Creates a cumulative probability of assignment plot for each	Error Checking

	NEWHYBRIDS result within a user-specified directory. Compliments <i>preCheckR</i> by allowing visually verification of Markov chain (non-) convergence.	and Diagnostics
<i>nh_plotR</i>	Plots the cumulative probability of assignment of a single NEWHYBRIDS result. Also allows the user to match plotting colours between analyses when NEWHYBRIDS reverses which population it designates Population 1 and 2.	Quantification and Analysis
<i>hybridPowerComp</i>	Evaluates the accuracy ³ and efficiency ³ with which NEWHYBRIDS assigns individuals of known hybrid class to the correct class across a range of minimum posterior probability thresholds from 0.50 to 0.99. Calculates the number of individuals wrongly assigned to hybrid genotype frequency classes over the total number of known pure individuals (type I error) ⁴ , and the proportion of individuals misclassified (type II error). The distribution of PofZ values for each genotype frequency class, as well as for all hybrid classes considered together is plotted. The effect of varying panel sizes on each of these variables is also evaluated. Plots are returned as .pdf and .jpg files, and all data frames constructed for plotting are exported.	Quantification and Analysis
<i>nh_accuracy_checkR</i>	Evaluates the accuracy with which NEWHYBRIDS assigns individuals of known hybrid class to the correct class for a single analysis at three	Quantification and Analysis

minimum posterior probability thresholds (PofZ \geq 0.05, 0.75 and 0.90).

This function is meant to compliment *hybridPowerComp*.

nh_panel_delta_plotR

Plots the genotype class assignment (class with max. PofZ) of individuals among panels of different size. Allows visualization of the stability of individual assignments to compliment the proportion of correct assignments returned by *hybridPowerComp*.

Quantification
and Analysis

nh_build_Example_Data

Writes example NEWHYBRIDS results to be evaluated with the function *hybridPowerComp*

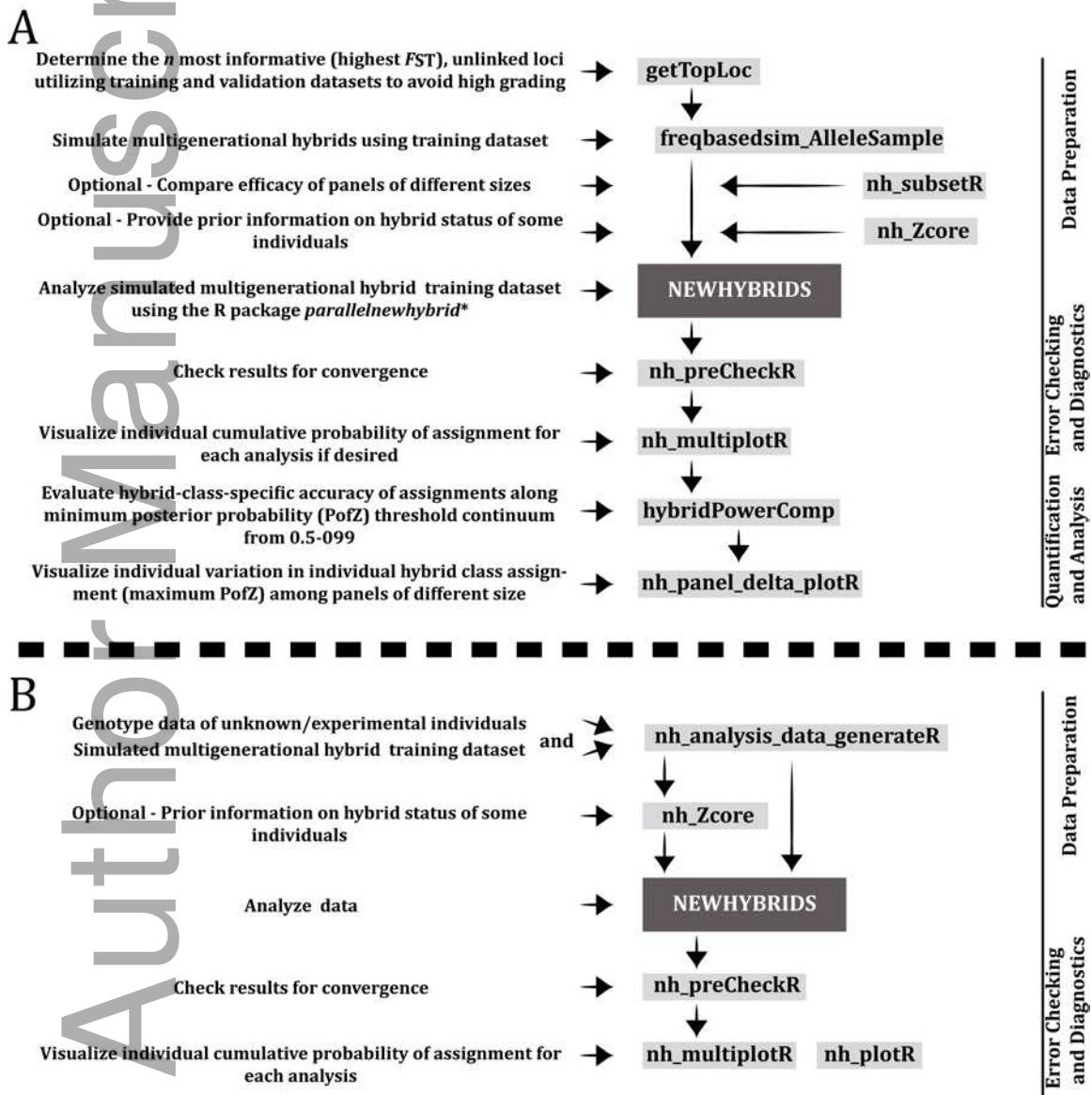
¹ Weir and Cockerham (1984)

² Anderson (2010)

³ Vaha and Primmer (2006)

⁴ Burgarella et al. (2009)

Figure 1. Schematic of the hybrid detection workflow and the associated functions (grey boxes) for: A, the development and quantification of the efficiency and accuracy of diagnostic panels of loci, and B the analysis of unknown/experimental data to detect hybrid individuals.



* (Wringe *et al.* 2017)

Figure 2. Plot of the efficiency of assignment for each of the six genotype frequency classes at critical posterior probability of assignment thresholds from 0.5 to 1.0 for diagnostic panels of various size. Each genotype frequency class is show in an individual facet, with abbreviations at its top. The solid coloured lines are the mean efficiency, and the dotted line the standard deviation of three independently simulated datasets, each analyzed in triplicate. Panel sizes and their corresponding colours are shown in the legend. The x-axis is the posterior probability of assignment threshold.

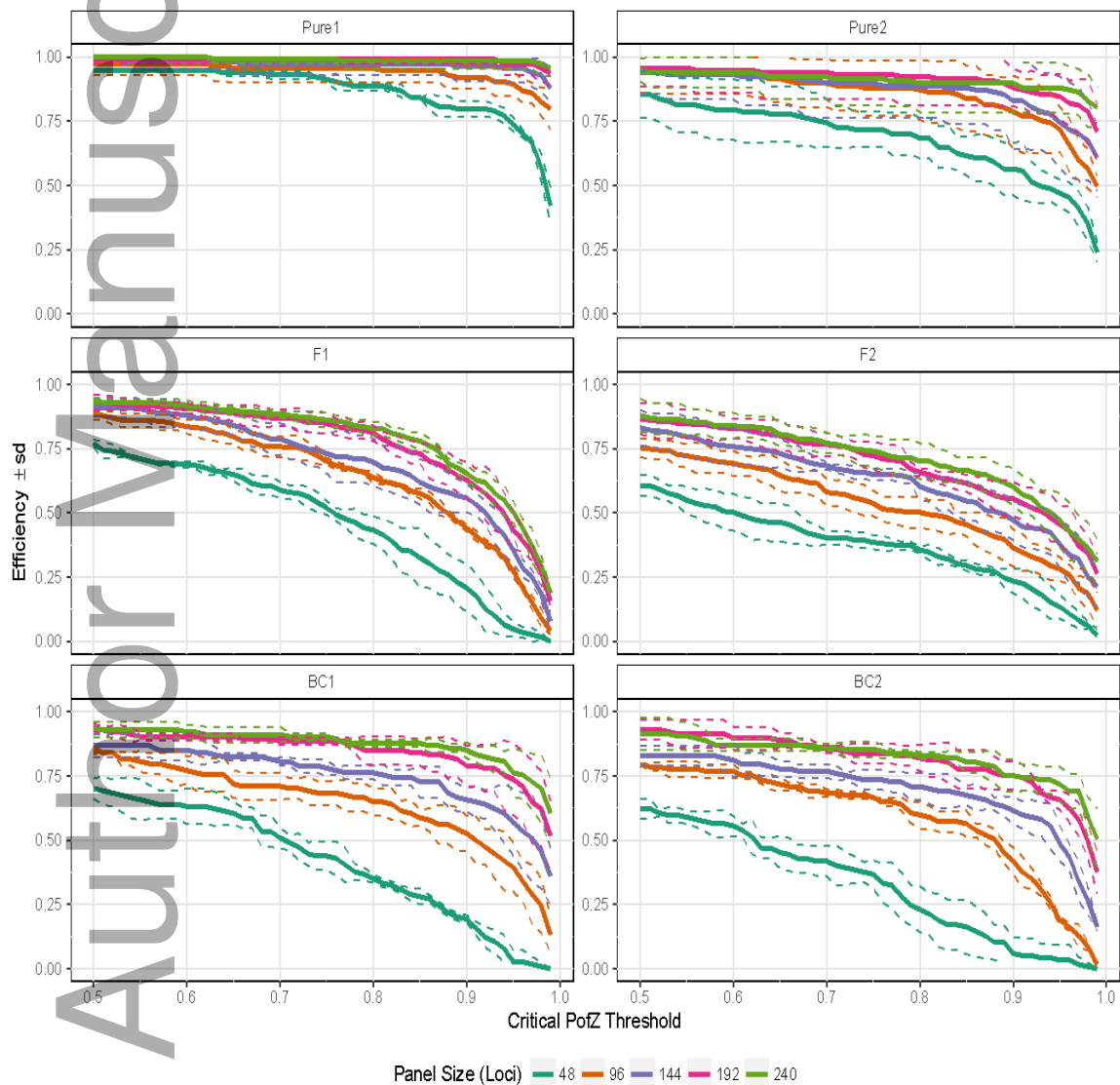


Figure 3. Plot of accuracy of assignment for each of the six genotype frequency classes for various panel sizes at critical posterior probability of assignment threshold values ranging from 0.5 to 1.0. Genotype frequency class abbreviations are as in Supplementary figure 1, and each class is displayed in a single facet. The solid coloured lines are the mean accuracy, and the dotted lines the standard deviation of three independently simulated datasets, each analyzed in triplicate. The panel sizes, and their representative colours are shown in the legend. The x-axis is the critical posterior probability of assignment threshold.

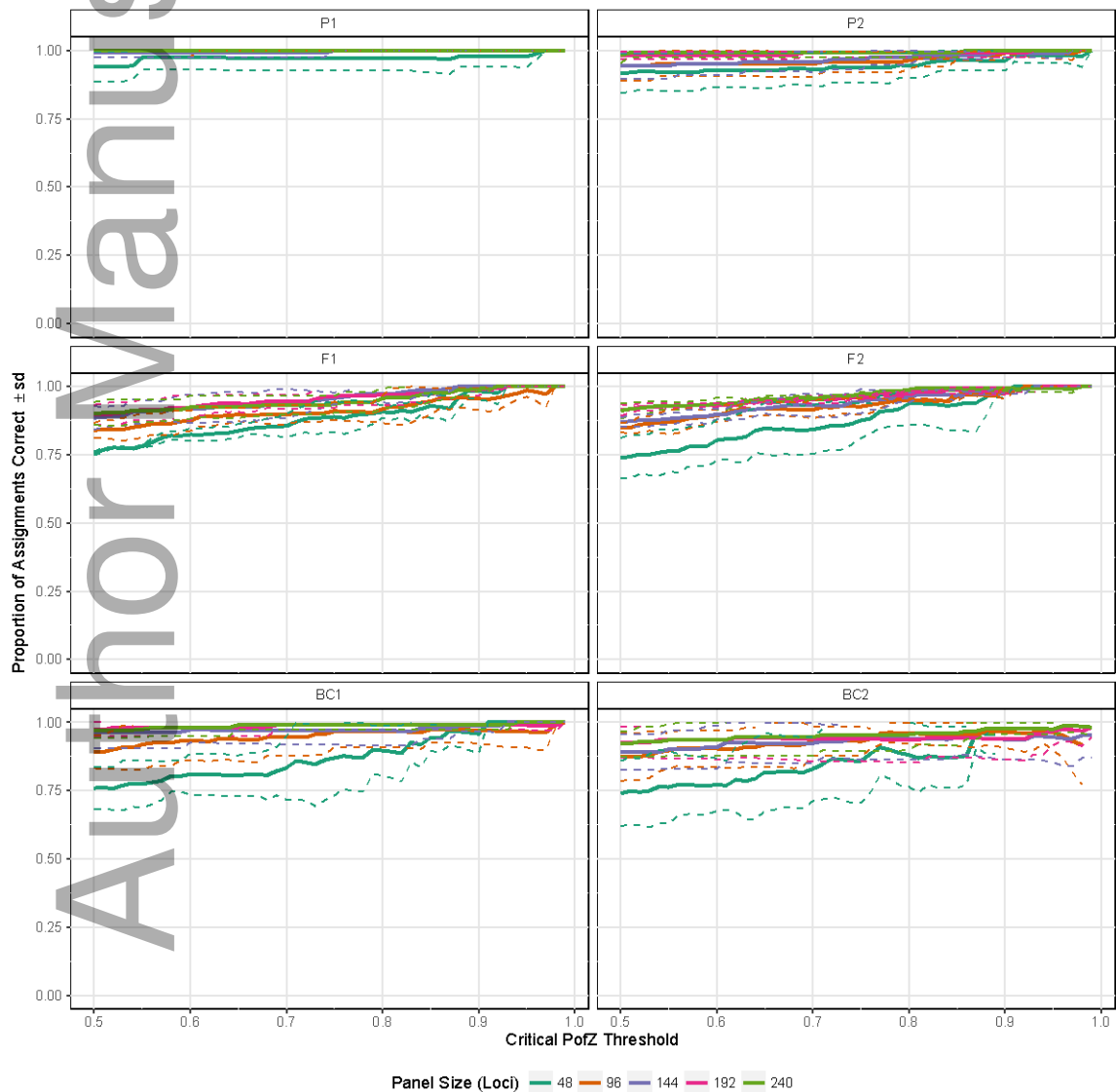


Table 1 – Functions included in the *hybriddetective* R package, a synopsis of their purpose, and which of the three major elements they are used in.

Function Name	Synopsis	Main Use
<i>getTopLoc</i>	Creates a panel comprised of the n (user-specified) most informative (based on highest loci-specific F_{ST} s), markers not in linkage disequilibrium. The function randomly assigns half the individuals in each of the two populations to be used to calculate loci-specific Weir and Cockerham's F_{ST} s ¹ , and returns the genotypes at the n loci of the other half to be used to test the efficacy of the panel to avoid high-grading bias ² .	Data Preparation
<i>freqbasedsim_GTFreq</i>	Creates simulated multi-generational (i.e. Pure 1, Pure 2, F ₁ , F ₂ , BC1, BC2) hybrids based on the allele frequencies in the two populations provided. The user can specify the number of individuals in each of the hybrid classes to be created.	Data Preparation
<i>freqbasedsim_AlleleSample</i>	Creates simulated multi-generational hybrids by randomly sampling, without replacement, two alleles per loci from a proportion of the individual genotypes provided. The user is able to specify the proportion of genotypes to sample, as well as the number of individuals of each hybrid class to create.	Data Preparation
<i>nh_analysis_generateR</i>	Merges a file composed of simulated hybrid genotypes with a file	Data

	containing the genotypes of unknown/experimental individuals to produce a file suitable to ascertain the hybrid class of the unknowns. The user is able to specify which hybrid classes from the simulated dataset to include in the output.	Preparation
<i>nh_analysis_simulateReferenceR</i>	Creates a simulated multi-generational hybrid reference dataset from user provided data, and then merges it with the genotypes of unknown/experimental individuals. This function will create a new simulated dataset each time it is run using the same simulation methodology as <i>freqbasedsim_GTFreq</i> .	Data Preparation
<i>nh_subsetR</i>	Removes subsets of desired loci from NEWHYBRIDS formatted files so that the efficacy of panels of various sizes can be assessed.	Data Preparation
<i>nh_Zscore</i>	Allows the user to assign known hybrid category designations to individuals in NEWHYBRIDS formatted files	Data Preparation
<i>nh_preCheckR</i>	Checks all NEWHYBRIDS results within a directory and flags those that show evidence that the Markov chain may have failed to converge. This is done by evaluating the proportion of known Pure Population 1 or 2 individuals in which the posterior probability of assignment to F ₂ exceeds a threshold. The user may specify both the proportion of individuals and the PofZ threshold.	Error Checking and Diagnostics
<i>nh_multplotR</i>	Creates a cumulative probability of assignment plot for each	Error Checking

	NEWHYBRIDS result within a user-specified directory. Compliments <i>preCheckR</i> by allowing visually verification of Markov chain (non-) convergence.	and Diagnostics
<i>nh_plotR</i>	Plots the cumulative probability of assignment of a single NEWHYBRIDS result. Also allows the user to match plotting colours between analyses when NEWHYBRIDS reverses which population it designates Population 1 and 2.	Quantification and Analysis
<i>hybridPowerComp</i>	Evaluates the accuracy ³ and efficiency ³ with which NEWHYBRIDS assigns individuals of known hybrid class to the correct class across a range of minimum posterior probability thresholds from 0.50 to 0.99. Calculates the number of individuals wrongly assigned to hybrid genotype frequency classes over the total number of known pure individuals (type I error) ⁴ , and the proportion of individuals misclassified (type II error). The distribution of PofZ values for each genotype frequency class, as well as for all hybrid classes considered together is plotted. The effect of varying panel sizes on each of these variables is also evaluated. Plots are returned as .pdf and .jpg files, and all data frames constructed for plotting are exported.	Quantification and Analysis
<i>nh_accuracy_checkR</i>	Evaluates the accuracy with which NEWHYBRIDS assigns individuals of known hybrid class to the correct class for a single analysis at three	Quantification and Analysis

minimum posterior probability thresholds (PofZ \geq 0.05, 0.75 and 0.90).

This function is meant to compliment *hybridPowerComp*.

nh_panel_delta_plotR

Plots the genotype class assignment (class with max. PofZ) of individuals among panels of different size. Allows visualization of the stability of individual assignments to compliment the proportion of correct assignments returned by *hybridPowerComp*.

Quantification
and Analysis

nh_build_Example_Data

Writes example NEWHYBRIDS results to be evaluated with the function *hybridPowerComp*

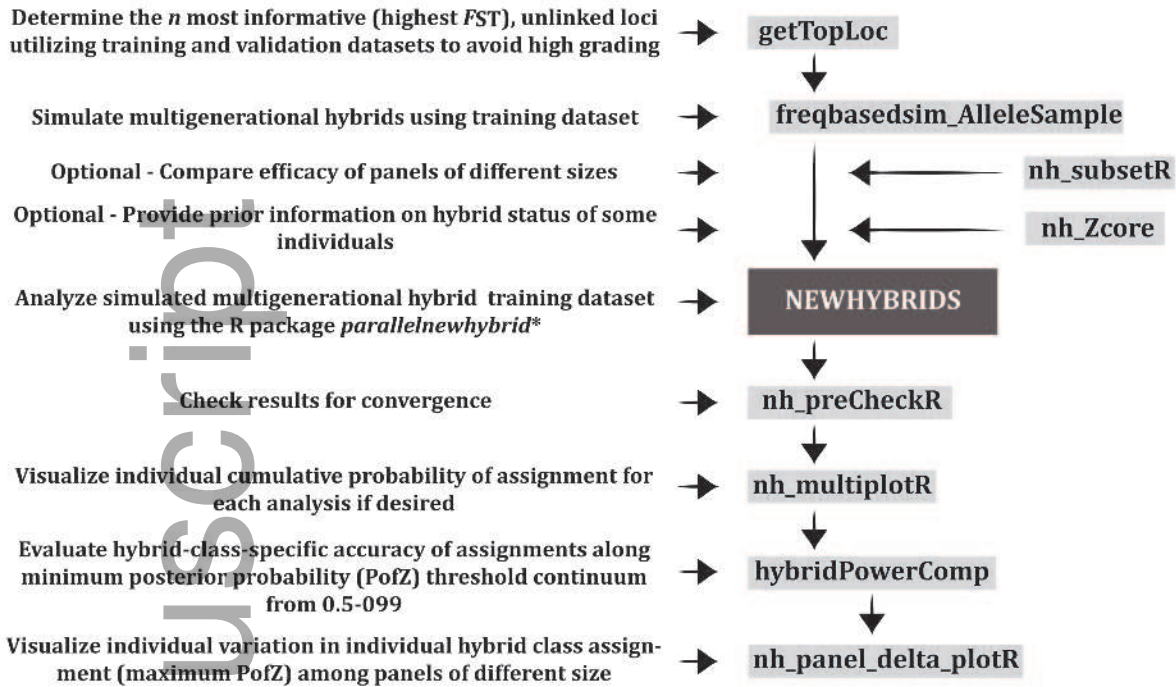
¹ Weir and Cockerham (1984)

² Anderson (2010)

³ Vaha and Primmer (2006)

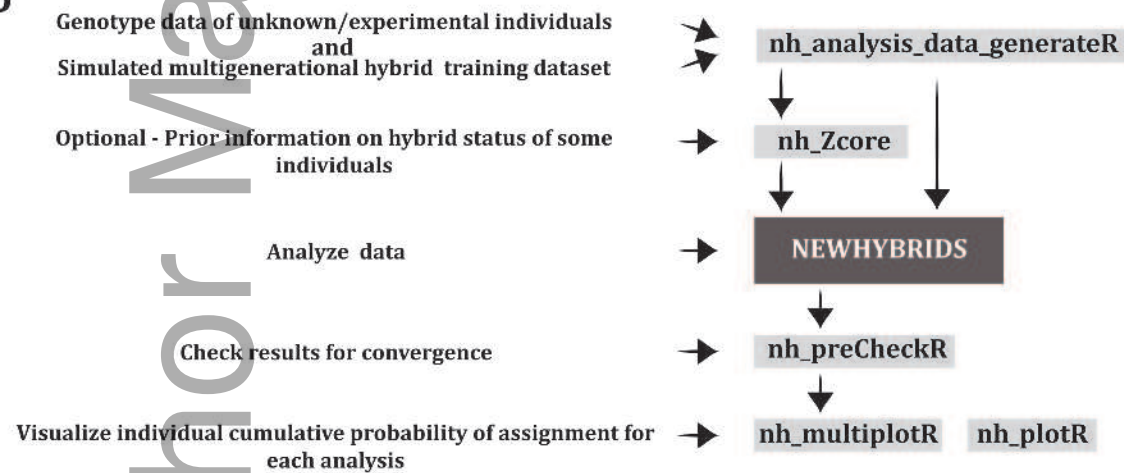
⁴ Burgarella et al. (2009)

A



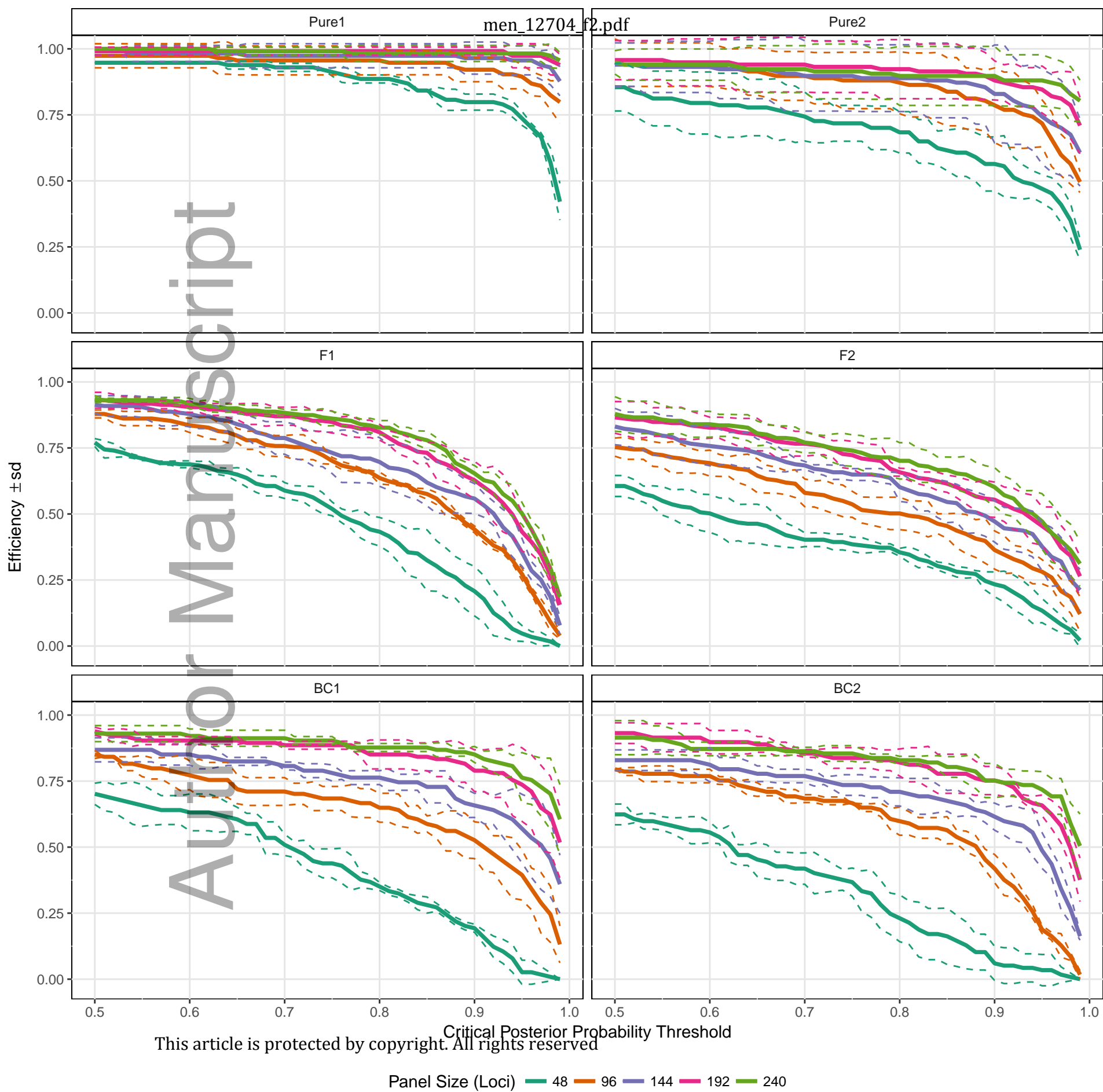
Data Preparation
Error Checking and Diagnostics
Quantification and Analysis

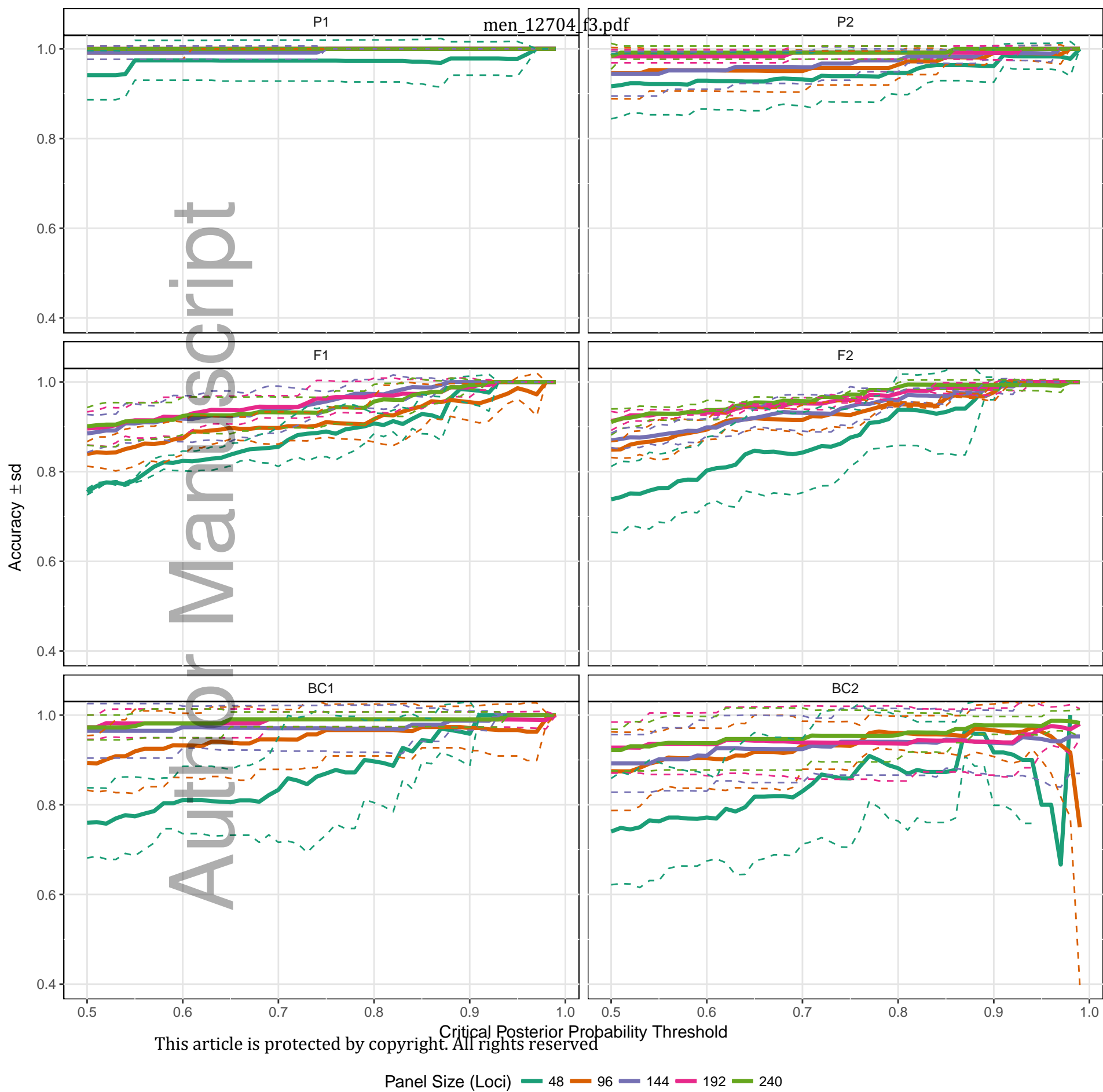
B



Data Preparation
Error Checking and Diagnostics

men_12704_f1.tif





This article is protected by copyright. All rights reserved

Panel Size (Loci) — 48 — 96 — 144 — 192 — 240