# Automated scoring of AFLPs using rawgeno v 2.0, a free R CRAN library

**3 authors**, including:

Nils Arrigo
University of Lausanne
**179** PUBLICATIONS   **1,806** CITATIONS

SEE PROFILE

Nadir Alvarez
Natural History Museum of Geneva
**384** PUBLICATIONS   **3,403** CITATIONS

SEE PROFILE

**Some of the authors of this publication are also working on these related projects:**

Project   ECOGEN - Ecosystem change and species persistence over time: a genome-based approach View project

Project   Selaginella View project

# Automated scoring of AFLPs using RawGeno v 2.0, a free R CRAN library

**Authors:** Arrigo Nils[1,2], Holderegger Rolf[3], Alvarez Nadir[2]

[1]Laboratory of Evolutionary Botany, Institute of Biology, University of Neuchâtel, 11 rue Emile-Argand, 2009 Neuchâtel, Switzerland

[2]Department of Ecology and Evolution, University of Lausanne – UNIL Sorge, Biophore Building, 1015 Lausanne, Switzerland

[3]WSL Swiss Federal Research Institute, Zürcherstrasse 111, 8903 Birmensdorf, Switzerland

## Abstract

Amplified Fragment Length Polymorphisms (AFLPs) are a cheap and efficient protocol for generating large sets of genetic markers. This technique has become increasingly used during the last decade in various fields of biology, including for instance population genomics, phylogeography and genome mapping. Here, we present RawGeno, an R library dedicated to the automated scoring of AFLPs (i.e. the coding of electropherogram signals into ready-to-use datasets). Our program includes a complete suite of tools for binning, editing, vizualizing and exporting results obtained from AFLP experiments. RawGeno can either be used with command lines and program analysis routines or through a user-friendly graphical user interface. We describe

27    the whole RawGeno pipeline along with recommendations for (i) setting the analysis of

28    electropherograms in combination with PeakScanner, a program freely distributed by Applied

29    Biosystems, (ii) performing quality checks, (iii) defining bins and proceeding to scoring, (iv)

30    filtering non-optimal bins and (v) exporting results in different formats.

31

32    **1. Introduction**

33    The amplified fragment length polymorphism (AFLP) technique is increasingly used in

34    phylogeographic and population genomics studies, particularly in non-model organisms for which

35    no prior DNA sequence information is available *(1)*. This relatively cheap technique is based on

36    complete endonuclease restriction digestion of total genomic DNA followed by selective PCR

37    amplification and electrophoresis of a subset of fragments, resulting in a unique, (theoretically)

38    reproducible fingerprint for each individual. Although the AFLP technique is able to generate a

39    large number of informative markers, the success of this method is compromised by different

40    factors *(2)*. For instance, manual scoring relies on visual inspection and subjective interpretation of

41    the electropherotic profiles during a time-consuming and tedious task. In the last decades, several

42    improvements in automatic scoring have been proposed and implemented in commercial software

43    [see *(3)* for a review]. However, until recently, no free open-source software was available to

44    process AFLP data from raw data to ready-to-use presence/absence binary matrices. Two years ago,

45    we developed RawGeno 1.0 *(4)*, a library performing automated binning, scoring and data mining

46    analyses under the R CRAN environment, based on outputs from the freely available

47    electropherogram-analyzing software PeakScanner (Applied Biosystems, Foster City, USA,

48    http://www.appliedbiosystems.com/peakscanner). Implementing RawGeno 1.0 solutions in a free

49    environment has provided an accessible and accurate solution to many users (with 415 downloads

50    from http://sourceforge.net/projects/rawgeno one year after release). Here, we present RawGeno

51    2.0, an updated version of this software, built around an optimized, less time-consuming algorithm,

52    implementing new features such as binning editing. We provide examples both in a user-friendly

53    and in a command-line interface. In order to allow users to customize queries (which may vary

54  depending on dataset quality/size), we further provide tips for setting analyses and improving output

55  robustness. All stages of a whole analysis are detailed according to the following five sections: (i)

56  importing features; (ii) quality check; (iii) binning and scoring algorithms; (iv) bin filtering; (v)

57  exporting options.

58

59  **2. Program Utilization**

60  *2.1. Overview of the analysis*

61  Analysing AFLP electropherograms is achieved using two programs: PeakScanner and

62  RawGeno. Whereas PeakScanner detects AFLP peaks along electropherograms and calculates their

63  intensity and size in base pairs (by relying on an internal size standard included in electrophoresis),

64  RawGeno proceeds to the binning and scoring of AFLP electropherograms.

65  RawGeno includes several filters to assess the quality of electropherograms and checks the

66  consistency of binning. In addition, several preliminary analyses are available to the user for

67  making biological inferences and/or remove outlier samples. Finally, several functions allow

68  exporting resulting data sets into properly formatted files for further analyses.

69

70  *2.2. Building up a scoring project*

71  In RawGeno 2.0 the AFLP scoring project should be organized according to the following

72  procedure.

73  a.  Create a folder (hereafter "project folder") from which the project will be managed.

74  b.  In this folder, add a sub-directory including all electropherogram files (*.fsa). Also add an R

75      shortcut (Windows users) to conveniently launch RawGeno scoring sessions. Right-clicking

76      on this R shortcut allows defining the default working directory of R by specifying it into the

77      "Start into" addressing field of the shortcut. Copy-pasting the project folder address into this

78      field will set-up the working directory of R accordingly.

79  c.  Create a text-tabulated table listing individuals included in the project (hereafter referred to

80      as "info table"). The info table is optional as the minimal RawGeno analysis can proceed

81    without it. However, RawGeno includes several functions relying on this table, for instance

82    to label individuals during preliminary analyses or facilitate the sorting and selection of

83    individuals (for example, according to populations or species) during the production of

84    exports. Therefore, the info table should include any additional relevant information that the

85    user would like to consider. It must contain at least the name of individuals (i.e. in a column

86    named "Tag") and any supplementary information in extra columns.

87

88    *2.3. Obtaining the raw data from \*.fsa files using PeakScanner*

89     The analysis of AFLPs starts by using PeakScanner in order to detect peaks along

90    electropherograms and to calculate their size. The procedure is highly automated, letting the user set

91    peak detection parameters and check the quality of electropherograms. The peak detection

92    parameters are set up using a so-called "Analysis Method", which is available from the graphical

93    interface in PeakScanner (menu "Resources/Manage Analysis Methods"). Typically, a proper peak

94    detection attempts to detect only peaks that are biologically relevant and exclude peaks only

95    reflecting technical background noise.

96     We advise to set the "Analysis Method" using the following guidelines.

97    a.  Prior to the detection of peaks *per se*, a light smoothing of electropherograms might be

98        desirable, in order to eliminate small secondary peaks due to technical background noise.

99    b.  The detection of peaks is achieved through a "sliding window" analysis that inspects

100       electropherograms locally. Within the inspected region, PeakScanner first creates a modelled

101       version of the electropherogram by fitting a polynomial curve to the data. Peaks are detected

102       according to this modelled signal, based on their absolute width. Therefore, the detection

103       sensitivity is adjusted by modifying the width of the sliding window (i.e. in terms of data

104       points, the smaller it is, the more sensitive the procedure becomes), the goodness of fit

105       reachable by the polynomial curve (again, increasing the polynomial degree of fitting

106       increases the detection sensitivity) and the minimal width above which a peak is recorded as

107       present. We advise to use the default parameters as a starting point, as they have been shown

4

108      to provide reliable results *(4, 5)*: set 15 points for the sliding window width, use a third

109      degree polynomial curve and consider peaks that at least have two measurement points of

110      half-width.

111    c.   Downstream to peak detection, PeakScanner filters peaks according to their absolute

112      fluorescence intensity, i.e. the peak height, measured in relative fluorescent units (rfu).

113      Visually checking electropherograms obtained from blank samples generally helps to adjust

114      the fluorescence threshold to the upper limit of the technical background noise. While some

115      applications might benefit from considering only peaks with a strong fluorescence (e.g.

116      greater than 150 rfu to provide conservative estimates for band presence statistics), most

117      users will prefer using a more permissive threshold at this stage and apply *a posteriori*

118      filtering strategies based on bin quality statistics *(6, 7)*. We advise to use 50 rfu as a minimal

119      fluorescence for considering individual AFLP peaks.

120    d.   Save the customized "Analysis Method" in order to use it during electropherogram analysis.

121    e.   Once the "Analysis Method" is set up, import the electropherograms (stored as *.fsa files)

122      into PeakScanner using the "Add Files" button.

123    f.   Define the size standard and the "Analysis Method" to be used for all individuals included in

124      the project (set this information for the first individual, then select the columns "Size

125      Standards" and "Analysis Method" and use the "ctrl+D" keyboard shortcut to apply these

126      settings to the remaining individuals).

127     The detection and sizing of peaks is processed using the "Analysis" menu, from the graphical

128 interface. Once achieved, electropherograms can be visualized and compared among individuals.

129 This might help identifying AFLP reactions that were not successful (e.g. individuals with a

130 systematically low fluorescence or showing abnormal peaks). Removing such individuals prior to

131 the RawGeno analysis will help to enhance the final quality of scoring.

132     The PeakScanner analysis ends with a simple export process in which the list of peaks detected

133 throughout the complete set of analyzed individuals is stored in a table. This is achieved using the

134 "Export/Export Combined Table" menu, producing a text-tabulated file containing the size, height,

135 area and width of all detected peaks (this can be checked using the "Edit Table Settings" menu).

136

137 ***2.4. From raw data to ready-to-use matrices using RawGeno***

138 *2.4.1. Installing RawGeno*

139 RawGeno is freely available from http://sourceforge.net/projects/rawgeno as a zip file. In

140 windows, the installation is achieved either using the graphical user interface of R (menu

141 "Packages/Install package(s) from local zip files") or the following command line in the R console:

142 `utils:::menuInstallLocal()`

143

144 Installing RawGeno with Linux requires decompressing RawGeno.zip into the library folder of

145 R. Using the shell command line, this is done as follows:

146 `sudo mv RawGeno.zip /usr/lib64/R/library/RawGeno.zip`

147 `cd /usr/lib64/R/library/RawGeno.zip`

148 `unzip RawGeno.zip`

149 `rm RawGeno.zip`

150

151 Finally, RawGeno requires the installation of two companion packages: vegan and tkrplot, that

152 both are available from usual R CRAN repositories (Linux users should see Note 1). Their

153 installation is achieved either using the graphical user interface of R (menu "Packages/Install

154 package(s) from CRAN") or with the following command lines (prompted into the R console):

155 `install.packages("vegan")`

156 `install.packages("tkrplot")`

157

158 *2.4.2. Importing PeakScanner results*

159 All following steps are performed in the R CRAN environment, using the R shortcut described

160 above (or ensuring that the correct working directory has been selected). Once the RawGeno

161 package has been installed, it can be called applying the following command line:

162    a.  Call RawGeno, vegan and tkrplot as libraries into R and launch the graphical user interface

163    ```
       require(RawGeno)
       ```

164    ```
       require(vegan)
       ```

165    ```
       require(tkrplot)
       ```

166    ```
       RawGeno()
       ```

167

168    b.  Importing the PeakScanner text-tabulated file in RawGeno can then be done using the

169        graphical user interface (menu "RawGeno/1. Files/Electroph./PeakScanner (*.txt)") or using

170        the following command lines:

171        Choose interactively the PeakScanner file

172    ```
       myfile=tk_choose.files(caption='Choose PeakScanner File')
       ```

173    ```
       OPENAFLP(myfile, pksc=T, dye="B")
       ```

174

175        Or explicitly specify the path of file of interest

176    ```
       mypath="C:/MyDocuments/MyPeakScannerFile.txt"
       ```

177    ```
       OPENAFLP(mypath, pksc=T, dye="B")
       ```

178

179    During importation, RawGeno handles a single dye color at a time, which is user-specified and

180    considers the "dye" parameter with the following values: "B" (blue; FAM), "G" (green; HEX), "Y"

181    (yellow; NED), "R" (red; ROX) or "O" (orange; LIZ). If electrophoresis was achieved using several

182    dyes simultaneously (e.g. multiplexing of PCR products), each dye must be analysed separately in

183    RawGeno. Datasets obtained from several dyes can be merged *a posteriori* in a final binary table

184    (see below).

185

186    *2.4.3. Quality check*

187    Because the detection of AFLP peaks is based on a defined threshold, it is not easy to handle

188    reactions showing electropherograms with varying intensities (see Note 2). When improperly

189  handled, such a situation leads to the inclusion of samples characterised by many false-absences in

190  the final dataset. Although the only way to correctly address this issue is a robust wet-lab protocol,

191  RawGeno still attempts limiting the influence of low quality AFLPs on binning and scoring, by

192  filtering individuals that were unsuccessful. Here, the variability in the number of peaks detected

193  per individual is used as a proxy of AFLP reactions quality. Empirical evidence shows that this

194  statistics is dependent of the specific dataset used and the biological organism studied. The lower

195  bound of this distribution most generally includes individuals with low AFLP intensities,

196  characterized by many AFLP peaks that remain undetected in the electropherograms. Because such

197  individuals usually represent a small fraction of the complete project, we advise removing them

198  from the dataset. The upper bound of the distribution can either reflect a biologically relevant signal

199  (e.g. hybridization) or a technical bias (e.g. contamination, odd PCR reaction). Such individuals

200  should be either discarded or identified as outliers for proper interpretation in further analyses.

201   RawGeno includes the two following options for checking quality of individuals:

202  a.  Filter individuals according to the number of detected peaks, by manually selecting

203     individuals that should be kept for further analysis or by using a dedicated device (Fig 1 a).

204     The command line version of this operation relies on percentiles and conserves individuals

205     within 5%-95% bounds of the detected peaks distribution.

206      Retrieve imported electropherograms

207     ```
data.electroph=AFLP$all.dat
```

208

209     Compute number of AFLP peaks per individual

210     ```
pk.per.smp=table(data.electroph$sample.id)
```

211

212     Compute 5% and 95% quantiles

213     ```
qtles=quantile(pk.per.smp, probs=c(0.05, 0.95))
```

214

215     Determine what samples can be kept

```
216        to.keep=which(pk.per.smp>=qtles[1] & pk.per.smp<=qtles[2])

217

218        Filter electropherograms, keep only retained individuals and update individuals indexing

219        accordingly

220        smp.ok=match(data.electroph$sample.id, to.keep)

221        data.clean=data.electroph[is.na(smp.ok)==F, ]

222        data.clean$sample.id=as.factor(data.clean$sample.id)

223        levels(data.clean$sample.id)=1:length(to.keep)

224        smp.names=AFLP$samples.names

225        smp.clean=smp.names[to.keep]

226        AFLP$all.dat=data.clean

227        AFLP$samples.names=smp.clean

228
```

229   b.  Check the quality of individuals, by taking into account their position in PCR plates (only

230       available from the graphical interface, Fig 1 b). This display helps to highlight systematic

231       biases having a technical origin (e.g. pipetting errors or thermocycler bias) and to identify

232       batches of individuals that were not successful.

233

234   *2.4.4. Binning*

235   Building a presence/absence matrix requires recognizing which AFLP peaks are homologous

236   across individuals. This procedure relies on the size of peaks along electropherograms and assumes

237   homology for peaks sharing identical sizes. Because peak sizes are determined empirically using

238   electrophoresis, measurements generally include technical variations preventing the observation of

239   strictly identical sizes across homologous AFLP peaks. Indeed, Holland *et al.* **(8)** reported

240   measurement variations ranging between 0 and 0.66 bp (with 0.08 bp in average) for replicated

241   AFLP peaks. Therefore, properly recording the signals of AFLP peaks asks for taking into account

242   size variations by defining size categories (i.e. "bins") into which the presence / absence of AFLP

243 peaks are recorded. Bins are characterised by their position along the electropherogram (i.e. the

244 average size of peaks they include) and their width (the size difference between the longest and

245 shortest peaks included in the bin). RawGeno uses a binning algorithm relying on the size of AFLP

246 peaks over all individuals included in the project and defines bins in a way that respects the two

247 following conditions (Fig 2 a).

248     a.  The first constrain is a maximal bin width. This prevents the definition of too large bins that

249        could lead to homoplasy (i.e. erroneously assigning non-homologous AFLP peaks within the

250        same bin). This limit is set using the MaxBin parameter. We advise to use MaxBin values

251        ranging between 1.5 and 2 bp (see Note 3). Using small values (i.e. MaxBin < 0.5 bp) should

252        be avoided as this generally causes oversplitting, a situation where the presence / absence of

253        homologous AFLP bands are coded using an exaggerated number of bins.

254     b.  The second constrain prevents the assignment of more than one peak from the same

255        individual within the same bin [i.e. "technical homoplasy", as defined in *(4)*]. If such a

256        situation occurs, RawGeno defines two separate bins in which the two peaks are assigned.

257        This constrain can be relaxed by increasing the "MinBin" parameter in order to include two

258        peaks of the same individual in the same bin (when not exceeding the "MinBin" value in size

259        difference). Such a relaxing might be desirable, for instance when artefactual peaks (i.e.

260        shoulder, stutter or secondary peaks bordering the authentic peak in a same individual) lead

261        to the definition of numerous extra-bins. In such a situation, artefactual peaks can cause the

262        local definition of extra bins into which homologous peaks can be inconsistently assigned.

263        We advise to use MinBin values ranging between 1 and 1.5 bp. The binning is launched

264        through the graphical interface (menu "RawGeno/2. Scoring") or using the following

265        command lines:

266        Proceed to binning

```
267        EXTRACTAFLP(all.dat=AFLP$all.dat,

268        samples.names=AFLP$samples.names, MAXBIN=2, MINBIN=1)
```

269

270    View results (that are assigned into a "data.binary" object)

271    `attributes(data.binary)`

272    `data.binary$data.binary`

273

274    Binning is an automated and straightforward analysis step that users might want to review

275    interactively. RawGeno includes a visualization device for manually editing the binning by adding,

276    removing or modifying the width and position of bins (this tool is only available from the graphical

277    interface, Fig 2 b). This device includes several help-to-decision statistics such as the average size

278    and the number of presences associated with each bin.

279

280    *2.4.5. Filtering*

281    Once defined, bins can be filtered according to their properties and/or quality. Note that such

282    filtering strategies require analyzing AFLP reactions with a consistent quality across individuals.

283    Filtering options are accessible from the scoring menu when using the graphical user interface

284    (menu "RawGeno/2. Scoring"). Command line users will set accordingly the EXTRACTAFLP

285    function.

286    Proceed to binning and filtering simultaneously

287    `EXTRACTAFLP(all.dat=AFLP$all.dat,`

288    `samples.names=AFLP$samples.names, MAXBIN=2, MINBIN=1, RMIN=100,`

289    `RMAX=500, cutRFU=50, who='B', thresh=95)`

290

291    In its current version, RawGeno includes three kinds of filters.

292    a.  The size filter restricts binning to a given portion of the electropherogram (RMIN and

293        RMAX parameters). We advise to limit the binning to peaks included in the range of the size

294        ladder, because their size is accurately interpolated by PeakScanner (in contrast to larger

295        peaks where the size is extrapolated). We recommend discarding peaks with small sizes (i.e.

296        smaller than 100 bp, RMIN=100) as they are more likely to be homoplasic *(9, 10)*. In

11

297       addition, large size peaks should as well be considered cautiously because their fluorescence

298       intensity might not always be consistent across individuals. Because this upper limit might

299       vary according to datasets, we advise to run preliminary analyses and check

300       electropherograms to confidently determine it.

301     b.   The second filter eliminates bins according to their average fluorescence (i.e. the cutRFU

302       parameter). This filter assumes that bins with a high average fluorescence retrieve a more

303       consistent signal than bins with a low fluorescence. The rationale for this strategy is the

304       following. The fluorescence of an AFLP fragment largely determines its detection probability

305       during the PeakScanner analysis of electropherograms. Therefore, fragments that

306       systematically produce low fluorescences are more likely to be erroneously recorded as

307       absent from electropherograms as they might pass the threshold in some reactions but not in

308       others just by chance. The computation starts by normalizing fluorescence intensities across

309       samples using the sum normalization method *(6)*, before computing the average fluorescence

310       of each bin. Hence, keep in mind that the cutRFU parameter applies on normalized values

311       and does not scale with fluorescence measures provided in PeakScanner. Setting this filter is

312       dataset-dependent and we recommend running several trials before producing a definitive

313       dataset (see Note 4). Refer to the works of Whitlock *et al.(6)* and Herrmann *et al. (7)* for

314       more sophisticated filtering R scripts based on fluorescence intensities. Bridging RawGeno

315       with these algorithms is achieved by exporting fluorescence results instead of a binary

316       matrix. Once binning is achieved, use the following command lines.

317       Retrieve raw fluorescence results, stored in a matrix corresponding to the usual binary

318       matrix

319
```
mat.rfu=t(data.binary$data.height.raw)
```

320

321       For normalised fluorescences (sum normalization), use instead

322
```
mat.rfu=t(data.binary$data.height)
```

323

324    Prepare for export and save as a text-tabulated file (refer to programs' documentation for

325    properly formatting files).

326    `mat.rfu[is.na(mat.rfu)==T]=0`

327    `write.table(mat.rfu, "MyFluorescenceFile.txt", quote=F,`

328    `sep="\t")`

329

330    c.  The reproducibility filter evaluates bin quality according to their robustness across AFLP

331    reactions, by relying on replicated samples. This filter assumes that replicated individuals

332    were selected randomly from the original dataset, so as to scan the genetic diversity at best

333    (see Notes 2 and 4). Keep in mind that RawGeno identifies replicated individuals using their

334    names. Replicates must be named using the original individual name plus a suffix letter. The

335    suffix is matched using the "who" parameter of the filtering algorithm. As an example,

336    "mysample.fsa" and "mysampleB.fsa" are a pair of original-replicated samples, being

337    identified with a "B" suffix (therefore, set who = "B" when filtering). For each bin, RawGeno

338    compares original to replicated individuals and calculates the percentage of original-

339    replicated pairs for which the AFLP signal is successfully reproduced. Bins where

340    reproducibility cannot reach a satisfactory rate (i.e. the "thresh" parameter, a user-defined

341    reproducibility percentage) are eliminated from the final dataset.

342

343    *2.4.6. Review of results*

344    RawGeno offers two displays for exploring scoring results (menu "RawGeno/3. Quality

345 Check/Samples Checking"). The binary matrix can be directly visualized using a heatmap, showing

346 individuals sorted according to their genetic relatedness. Alternatively, individuals can be examined

347 with a principal coordinates analysis. Both displays allow plotting quality statistics or external

348 information (i.e. picked from the "info table" cited above) against the AFLP results. These displays

349 are only available from the graphical user interface (Fig 3 a, b).

350

351 *2.4.7. Exporting files*

352 RawGeno includes functions for producing binary tables and standard exports for Arlequin,

353 Hickory, Popgen, AFLPsurv, STRUCTURE 2.2, Mltr, Spagedi, Dfdist, Treecon, Baps, PAUP,

354 Structurama, MrBayes and NewHybrids *(11)*. Furthermore, these exports can be sliced according to

355 information provided in the info table and produce ad-hoc subsets. These functions are accessible

356 from the graphical interface ("RawGeno/4. Save") or using the following command lines.

357 Retrieve the binary matrix from RawGeno and import the info table

358 `mat01=t(data.binary$data.binary)`

359 `matinfo=read.delim("MyInfoTable.txt", header=T)`

360

361 Cross-reference the AFLP results to the info table

362 `popsA = row.names(mat01)`

363 `popsB = matinfo$Tag`

364 `mat01= mat01[match(intersect(popsA, popsB), popsA), ]`

365 `matinfo = matinfo[match(intersect(popsA, popsB), popsB), ]`

366

367 Remove monomorphic bins

368 `mat01=mat01[ , colSums(mat01)>0 & colSums(mat01)<nrow(mat01)]`

369

370 Produce the required outputs, e.g. for STRUCTURE 2.2. (refer to the library documentation for

371 further details regarding exporting functions).

372 `Structure.popsD(mat01, pops=matinfo$MyPopsColumn, path=getwd(),`

373 `name="MyStructure2.2File.txt")`

374

375 Users willing to analyse AFLP signals as codominant markers *(12)* should use command lines

376 described above to export fluorescence data (a proxy of allele copy number in genomes) associated

377 with binary matrices.

378

14

379 *2.4.8. Handling data from previously scored projects*

380     Users willing to merge, visualize and/or produce exports from datasets that were already scored

381 can import binary tables within RawGeno using the "RawGeno/1. Files/Import" menu. From the

382 command line, such an operation is done as follows.

383     Select files to merge

384 `list.merge=tk_choose.files(caption='Choose Files to Merge')`

385

386     Or specify a directory in which the binary matrices are stored.

387 `mypath="C:/MyDocuments/MyBinaryMatricesDirectory"`

388 `list.merge=dir(mypath, pattern='.txt')`

389

390     Proceed to merging

391 `MERGING(transpose = "indRows", exclude = T, replacewith = NA)`

392

393     The transpose parameter states whether the binary matrices store individuals as lines

394 ("indRows") or as columns ("indColumns"), the exclude parameter defines whether individuals that

395 are not shared by all matrices will be removed from the final merged dataset (exclude = "T"). If

396 kept (exclude = "F"), individuals with missing AFLP genotypes will be completed using NA values

397 (replacewith = NA) when no data is available. Note that the merged matrix is stored into a

398 "mergedTable" object. Visualization and exports can be performed using the graphical user

399 interface, as described above. Command lines for exporting merged matrices are given below.

400 `mat01=mergedTable`

401 `matinfo=read.table("MyInfoTable.txt", header=T)`

402 `popsA = row.names(mat01)`

403 `popsB = matinfo$Tag`

404 `mat01= mat01[match(intersect(popsA, popsB), popsA), ]`

405 `matinfo = matinfo[match(intersect(popsA, popsB), popsB), ]`

406 `mat01=mat01[ , colSums(mat01)&colSums(mat01)<nrow(mat01)]`

```
407    Structure.popsD(mat01, pops=matinfo$MyPopsColumn, path=getwd(),
408    name="MyStructure2.2File.txt")
409
```

**3. Conclusions and perspectives**

We present here a complete suite of tools to automate the scoring of AFLP datasets using free software applications Our program proposes an integrated solution to manage all the components of the analysis pipeline. Accordingly, samples are checked by removing non-satisfactory electropherograms at the very beginning of the analyses and AFLP genotypes are associated with user-specified information while producing ad-hoc exports. Bins are managed using an automated algorithm and can be edited manually using a dedicated graphical user interface.

As a next milestone, we plan to develop RawGeno into two complementary directions: incorporating the handling of electropherograms (which is now is part of PeakScanner) and the enhancement of bin filtering possibilities. Indeed, the RawGeno version currently under development already integrates functions for detecting and calculating the size of AFLP peaks along electropherograms. Finally, achieving connections with the R scripts of Herrmann *et al.* **(6)** and Whitlock *et al.* **(7)** is another way to improve RawGeno.

**4. Notes**

*Note 1*

Linux users might need to run R as "sudo" users to properly install companion packages (i.e.

vegan and tkrplot). In addition, troubles might arise because R libraries are downloaded as source

code and compiled locally before being installed. This requires that all compilers needed by R (such

as gc, gcc, gcc-fortran and others) have been installed locally, before attempting the installation of

external R packages. In OpenSUSE, the necessary compilers can be obtained using YaST2 (into the

rpm groups dedicated to development tools). Ubuntu users are more fortunate because Synaptic

Manager can install ready-to-use R libraries in addition to usual compilers (refer to http://cran.r-

project.org/bin/linux/ubuntu/README for further details regarding repository addresses).


*Note 2*

Whereas manual scoring allows permanent but subjective adjusting of the criteria defining

whether an AFLP peak should be recorded as present or absent, automated peak detection

algorithms apply uniform fluorescence thresholds. This requires datasets showing low AFLP quality

variation, because fluorescence differences between samples will be reflected in the final binary

matrix. For instance, we encountered problematic situations with AFLP reactions showing different

fluorescence offsets among PCR plates. These "plate effects" can be highlighted with principal

coordinate analyses, as samples are clustered according to PCR plates (i.e. due to plate-specific

losses of AFLP bands). These situations are especially difficult to handle without repeating wet-lab

experiments. Solving this problem computationally remains difficult and asks for setting sample (or

plate)-specific detection sensitivities, which is beyond PeakScanner possibilities. The present

version of RawGeno conservatively proposes the removal of samples with unusual AFLP profiles,

by relying on the distribution of the number of AFLP peaks. Future developments of RawGeno will

attempt analysing electropherograms directly. The problem highlights the crucial importance of

using robust and standardized lab protocols. In the following, we list several tips helping to retrieve

consistent results from AFLP reactions, when using automated scoring [also consult Bonin *et al.*

451 *(13)* and Gugerli *et al. (14)*].

    a.   Randomize reactions on PCR plates, in order to properly discriminate technical bias from

453        biological signals.

454     b.   Standardize all reaction steps: adjust DNA concentrations after spectrometer quantification,

455        limit impacts of pipetting errors by preparing reaction mixes in batches, run critical reactions

456        in uniform conditions (for instance perform restriction steps in an incubator rather than in a

457        thermocycler) and, importantly, run PCRs on the same thermocycler.

458     c.   Optimize signal detection during genotyping analysis, for instance by increasing the injection

459        time of automated sequencers at the beginning of electrophoresis, and prefer primer pairs

460        showing strong and consistent amplifications.

461

462 *Note 3*

463     Based on empirical case-study datasets (available from the authors upon request), we provide

464 help-to-decision statistics aimed at guiding users in setting their RawGeno analysis. We analysed 17

465 AFLP primer datasets: *Aegilops geniculata*, two primer pairs *(5)*; *Arum* spp., two primer pairs

466 (Espindola *et al.* unpubl. data); *Baldellia* spp., two primer pairs (Arrigo *et al.*, unpubl. data);

467 *Bupleurum ranunculoides*, three primer pairs (Labhardt *et al.*, unpubl. data); *Deschampsia litoralis*

468 and *Deschampsia caespitosa*, each with three primer pairs *(15)*; *Peucedanum ostruthium*, two

469 primer pairs (Borer *et al.*, unpubl. data). These studies included between 87 and 509 individuals

470 (mean: 255) of either intra- and interspecific sampling; 4% to 51% (mean: 23%) of the individuals

471 were replicated. We varied binning parameters (i.e. MinBin and MaxBin) and measured their effects

472 on the width of bins, the datasets' polymorphism (i.e. the proportion of bins with presence

473 frequencies ranging between 5% and 95%) and the bin reproducibility (i.e. the proportion of

474 replicated sample pairs over which the focal bin is successfully reproduced; this measure is

475 independent of the total number of bins in the dataset and is therefore suitable in the context of

476 binning optimization).

477    MinBin and MaxBin acted as boundaries on the width of bins. They determined how accurately

478 an AFLP signal was reflected into the final presence/absence matrix. Both parameters were

479 explored for values ranging between 0.1 bp and 5 bp, with 0.1 bp increments, therefore totalizing

480 1176 "binning trials" per dataset (Fig 4. a, b).

481     Using exaggeratedly small MaxBin values forces the binning algorithm to define narrow bins.

482 This situation leads to "oversplitting", a bias where AFLP signals are coded into more bins than

483 needed. In this case, an AFLP locus is coded using several adjacent bins appearing as inconsistent

484 when considered independently. Our results showed oversplitting evidence for bin widths below

485 0.5 bp, with decreased bins' polymorphism and reproducibility (Fig 4. c, d). This situation should be

486 avoided, and we recommend using values larger than 0.5 bp for MaxBin (in contrast, the MinBin

487 parameter had little effects on oversplitting).

488     Using exaggeratedly large MinBin and MaxBin values introduces "technical homoplasy" [as

489 defined in *(4)*], a bias where AFLP signals are coded using less bins than required. Although

490 merging artefactual secondary peaks with authentic peaks is desirable, a process that increases the

491 consistency of binning (see above), exaggerated merging tends to artificially increase similarity

492 between unrelated samples and has immediate effects on AFLP polymorphism and reproducibility.

493 Our results showed that technical homoplasy was reflected by a decrease in bin polymorphism,

494 when MinBin and MaxBin both exceeded 2 bp (i.e. corresponding to bin widths larger than 0.8 bp).

495 On the other hand, reproducibility increased along with technical homoplasy due to the addition of

496 some level of artefactual similarity among samples.

497     From these results, we suggest to screen binning parameters by considering bin polymorphism as

498 a main optimization criterion (Fig 4. a, c). Reproducibility statistics should not be considered for

499 binning optimization because of their inability to detect technical homoplasy (Fig 4. b, d). See

500 Holland *et al. (7)* for further considerations about binning optimization.

501

502     ***Note 4***

503     In its current version, RawGeno includes three filters that can be applied after binning has been

504 achieved. These are a size (i.e. the region of electropherograms over which the analysis must be

carried out), a fluorescence and a reproducibility filters. We applied filters to the 17 datasets

explored during binning optimizations (see Note 3). The size filter was not tested and all datasets

were analysed for bins ranging between 100 and 400 bp (i.e. corresponding to the electrophoresis

region where size interpolation is accurate). The two remaining filters were explored starting from

binned datasets that maximized polymorphism.

The first filter considers bins' quality to vary according to average fluorescence. Indeed, bins

with an average fluorescence close to the detection threshold used in PeakScanner are more likely

to reflect inconsistent signals (e.g. technical false absences) than bins with strong average

fluorescence (see above). The filter removes bins according to a user-defined lower fluorescence

limit. Keep in mind that this limit applies on normalized values. We tested limits ranging between 0

(i.e. no filtering) and 400 rfu (with 10 rfu increments) and measured polymorphism and

reproducibility changes caused by this filtering (Fig 5. a, b). All but one dataset gained in

polymorphism and reproducibility when filtering was optimized. However, while filtering increased

dataset reproducibility, it drastically decreased polymorphism when set up improperly. Its use hence

requires dataset-depending optimization to limit information reduction in datasets. We recommend

testing fluorescence thresholds below 200 rfu.

The reproducibility filter assesses the robustness of AFLP signals and removes bins that are not

satisfactorily reproducible. We tested filtering using bin reproducibility limits ranging between 0

and 100% (with 5% increments) and measured polymorphism and global error rate of datasets

(Fig 5. c, d). We present error rates *(13)* instead of bin reproducibility here because error rates are

more commonly reported in AFLP studies. Filtering effectively reduced the global error rate along

with the removal of non-reproducible bins from datasets. As expected, filtering also decreased the

polymorphism of datasets. Nevertheless, this filter generally removed numerous bins, and we

clearly propose avoiding the use of reproducibility values larger than 95%. Our experience shows

that using 85% as a threshold provides satisfactory results by balancing the error rates of datasets

(i.e. 5% in average, which we consider reasonable) with polymorphism. A more stringent filtering

can be applied (e.g. in studies where peak reproducibility is especially relevant such as in genome

532    scans), but it might affect the information content of datasets and thus requires the inclusion of

533    additional primer pairs to achieve a reasonable level of polymorphism.

**6. References**

1.  Vos P, Hogers R, Bleeker M, Reijans M, van de Lee T, Hornes M, Frijters A, Pot J, Peleman J, Kuiper M, Zabeau M (1995) AFLP: a new technique for DNA fingerprinting. Nucleic Acids Res 23:4407–4414.

2.  Pompanon F, Bonin A, Bellemain E, Taberlet P (2005) Genotyping errors: causes, consequences and solutions. Nat Rev Genet 6:847–859.

3.  Meudt HM, Clarke AC (2007) Almost forgotten or latest practice? AFLP applications, analyses and advances. Trends Plant Sci 12(3):106-117.

4.  Arrigo N, Tuszynski JW, Ehrich D, Gerdes T, Alvarez N (2009) Evaluating the impact of scoring parameters on the structure of intra-specific genetic variation using RawGeno, an R package for automating AFLP scoring. BMC Bioinformatics doi:10.1186/1471-2105-10-33.

5.  Arrigo N, Felber F, Parisod C, Buerki S, Alvarez N, David J, Guadagnuolo R (in press) Origin and expansion of the allotetraploid *Aegilops geniculata*, a wild relative of wheat. New Phytol.

6.  Whitlock R, Hipperson H, Mannarelli M, Butlin RK, Burke T (2008) An objective, rapid and reproducible method for scoring AFLP peak-height data that minimizes genotyping error. Mol Ecol Res 8:725-735.

7.  Herrmann D, Poncet BN, Manel S, Rioux D, Gielly L, Taberlet P, Gugerli F (2010) Selection criteria for scoring amplified fragment length polymorphisms (AFLPs) positively affect the reliability of population genetic parameter estimates. Genome 53(4):302-310.

8.  Holland BR, Clarke AC, Meudt HM (2008) Optimizing automated AFLP scoring parameters to improve phylogenetic resolution. Syst Biol 57(3):347-366.

9.  Vekemans X, Beauwens T, Lemaire M, Roldan-Ruiz I (2002) Data from amplified fragment length polymorphism (AFLP) markers show indication of size homoplasy and of a relationship between degree of homoplasy and fragment size. Mol Ecol 11(1):139-151.

10. Paris M, Bonnes B, Ficetola GF, Poncet BN, Després L (2010) Amplified fragment length

569       homoplasy: in silico analysis for model and non-model species. BMC Genomics

570       doi:10.1186/1471-2164-11-287.

571   11. Ehrich D (2006) AFLPdat: a collection of R functions for convenient handling of AFLP

572       data. Mol Ecol Notes 6: 603-604.

573   12. Gort G, van Eeuwijk F (2010) Codominant scoring of AFLP in association panels. Theor

574       Appl Genet doi:10.1007/s00122-010-1313-x.

575   13. Bonin A, Bellemain E, Eidesen PB, Pompanon F, Brochmann C, Taberlet P (2004) How to

576       track and assess genotyping errors in population genetics studies. Mol Ecol 13:3261–3273.

577   14. Gugerli F, Englisch T, Niklfeld H, Tribsch A, Mirek Z, Ronikier M, Zimmermann NE,

578       Holderegger R, Taberlet P, IntraBioDiv Consortium (2008) Relationships among levels of

579       biodiversity and the relevance of intraspecific diversity in conservation - a project synopsis.

580       Perspect Plant Ecol Evol Syst 10: 259-281.

581   15. Peintiger M, Arrigo N, Brodbeck S, Koller A, Imsand M Holderegger R (2010) Genetische

582       und morphologische Differenzierung der endemischen Grasart *Deschampsia littoralis*

583       (Gaudin) Reut. – Wie verschieden sind die Population am Bodensee und am Lac de Joux im

584       Vergleich zu *D. cespitosa* (L.) P. Beauv.? Vogelwarte Radolfszell, University of Neuchâtel,

585       WSL.

**7. Figure captions**

587     Fig 1. Checking quality of AFLP reactions. (a) Interactive device for filtering samples, before

588     performing binning. Samples are sorted according to the number of AFLP peaks successfully

589     amplified. Summary statistics (i.e. 5%, 50% and 95% quantiles) are provided to help users selecting

590     samples to be included in further analyses. (b) Visualization of 96-well PCR plates. Four quality

591     statistics are available in RawGeno (number of AFLP peaks per sample, mean and variance of

592     fluorescence intensity and outlier detection index; see RawGeno manual). Boxplots provide

593     summary statistics for rows and columns of PCR plates, respectively.

594

595     Fig 2. Binning algorithm implemented in RawGeno (from *(4)*). (a) Left panel: main steps

596     followed by the algorithm to define bins; right panel: illustration of binning with two samples (S1

597     and S2); the bin widths (i.e. the difference in size between the longest and the shortest amplicons

598     included in the considered bin) and the technical homoplasy rates (i.e. HR, the mean number of

599     peaks belonging to the same sample that are included in a same bin) are indicated. (b) Binning

600     edition device, allowing users to translate, resize, add or remove bins interactively, starting from

601     bins that were initially defined by the automated algorithm. Summary statistics (i.e. the average size

602     of bins and the number of AFLP peaks present per bin) are provided as editing guidelines.

603

604     Fig 3. Reviewing results. RawGeno includes basic visualization devices for performing

605     preliminary data mining. Specifically, results can be reviewed using (a) heatmaps of the binary

606     matrix, where samples are sorted according to their genetic similarity and (b) principal coordinates

607     analysis of the corresponding matrix. Both devices can compare AFLP results with either quality

608     statistics (i.e. number of AFLP peaks per sample, mean and variance in fluorescence intensity and

609     outlier detection index) or external information provided by users (e.g. the population from where

610     samples were collected). In addition, both devices are handled through a graphical user interface for

611     sorting and selecting samples to be visualized.

612

613     Fig 4. Binning parameters. We analysed 17 different datasets (see Note 3) in RawGeno, by

614     varying binning parameters that have an effect on the width of bins (i.e. MinBin and MaxBin) and

615     measured the associated effects on bin width, polymorphism and reproducibility. All results were

616     corrected in order to normalize statistics among datasets (i.e. using non-linear mixed effect models,

617     with the dataset origin considered as a covariable). Upper row: triangle plots displaying (a) bin

618     polymorphism and (b) reproducibility averaged for the 17 datasets as a function of MinBin and

619     MaxBin parameters. Lower row: scatterplots of corresponding (c) bin polymorphism and (d) bin

620     reproducibility statistics, displayed according to the median bin width of binning trials. 5% and 95%

621     confidence intervals (dashed lines) and the average (continuous line) are displayed. (a) to (d) Trials

622     optimizing bin polymorphism are represented as dots for each of the 17 datasets (boxplots indicate

623     bin widths associated to optimized trials).

624

625     Fig 5. Filtering parameters. We filtered datasets (using optimal binning parameters as defined in

626     Note 3) by increasing bin fluorescence and reproducibility thresholds. Median bin polymorphism

627     and error rates were measured for the final datasets (dataset-specific offsets were not corrected).

628     Upper row: bins are filtered according to their average fluorescence. Results are displayed as a

629     comparison with non-filtered trials, with "D bin polymorphism" and "D bin reproducibility" being

630     the difference between filtered and non-filtered datasets. Therefore, values greater than zero reflect

631     cases where filtering increased the median (a) polymorphism or (b) reproducibility of bins in the

632     final dataset. Conversely, values lower than zero indicate decreases in these statistics associated

633     with filtering. Lower row: filtered bins according to their reproducibility and consequences on

634     median bin (c) polymorphism and (d) error rate of the final datasets. Results of each dataset (grey

635     lines) are displayed along with summary statistics (dashed lines: 5% and 95% confidence intervals;

636     continuous line: average). (a) to (d) Trials optimizing bin polymorphism after filtering are presented

637     as dots for each of the 17 datasets (boxplots indicate fluorescence threshold associated to optimized

638     filtering).

(a) Clean Samples — AFLP peaks per individual. 375 individuals kept, 97% sampling kept. Conserved for further analysis (green region between 5% and 95%). Reference values: 195, 77.

(b) Heatmap of values across columns A–H and rows 1–12 with marginal boxplots.

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 12 | 109 | 117 | 82 | 83 | 38 | 78 | 94 | 50 |
| 11 | 97 | 107 | 84 | 93 | 105 | 96 | 111 | 92 |
| 10 | 49 | 104 | 95 | 107 | 66 | 104 | 65 | 55 |
| 9 | 104 | 101 | 53 | 101 | 91 | 96 | 98 | 78 |
| 8 | 58 | 74 | 115 | 79 | 109 | 98 | 49 | 49 |
| 7 | 86 | 101 | 85 | 57 | 92 | 82 | 95 | 64 |
| 6 | 83 | 94 | 88 | 39 | 78 | 87 | 90 | 58 |
| 5 | 87 | 105 | 85 | 61 | 94 | 84 | 105 | 86 |
| 4 | 113 | 95 | 105 | 98 | 104 | 90 | 60 | 64 |
| 3 | 61 | 86 | 104 | 97 | 70 | 89 | 92 | 97 |
| 2 | 112 | 105 | 81 | 83 | 72 | 82 | 78 | 73 |
| 1 | 125 | 106 | 62 | 30 | 60 | 41 | 75 | 1 |

(a) (b)

(a) Median bin polymorphism

(b) Median bin reproducibility

(c)

(d)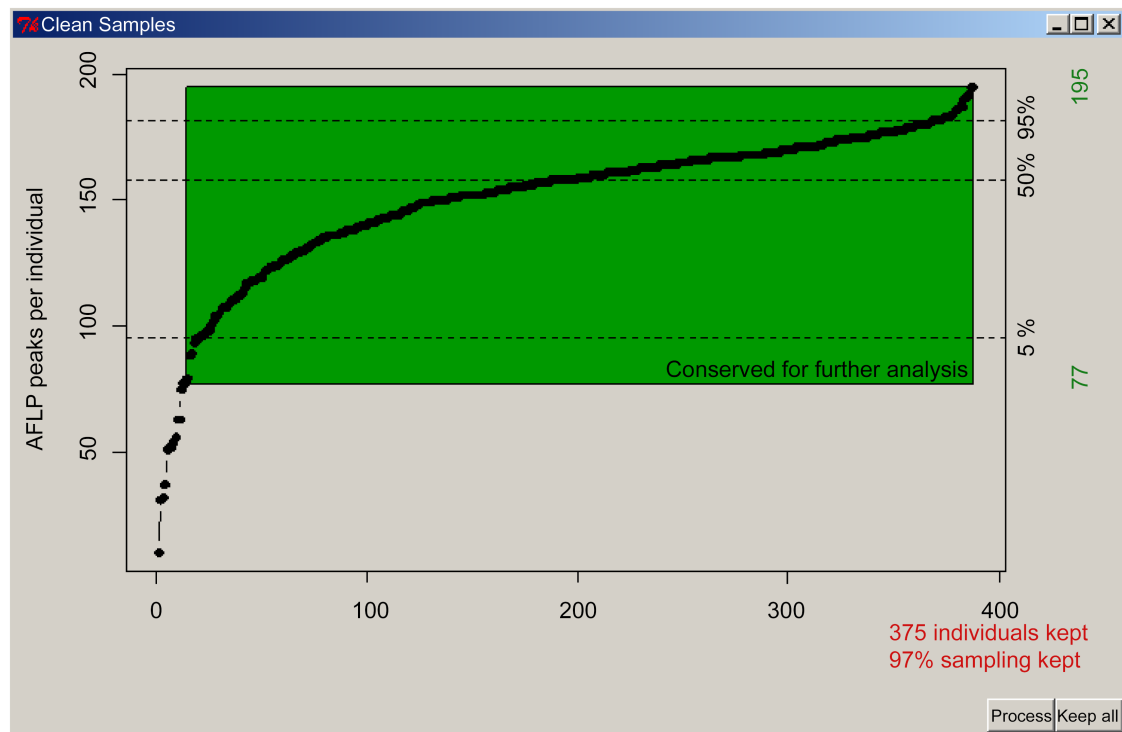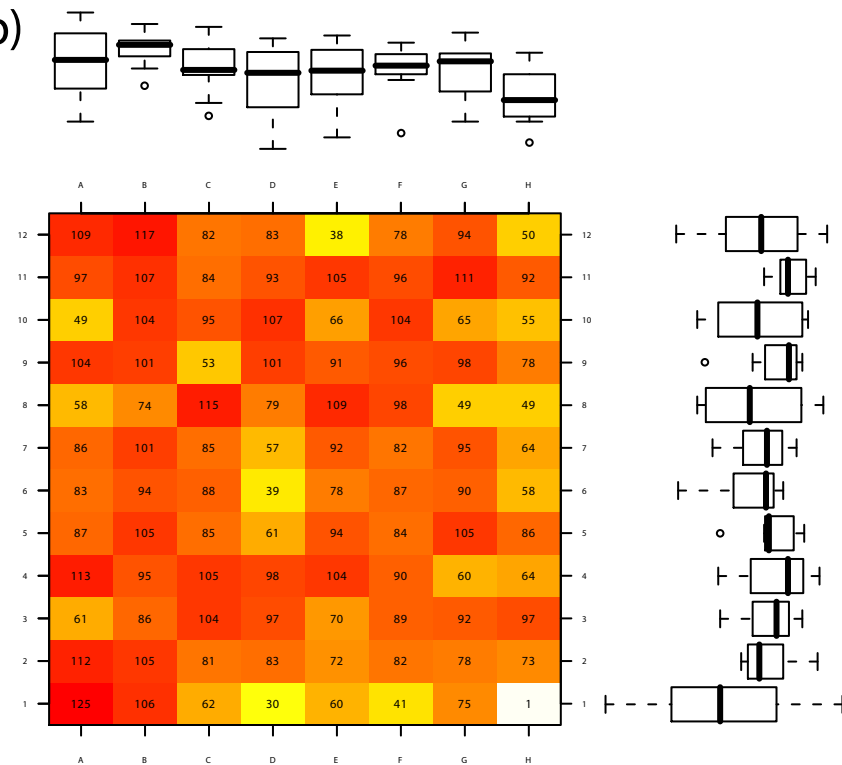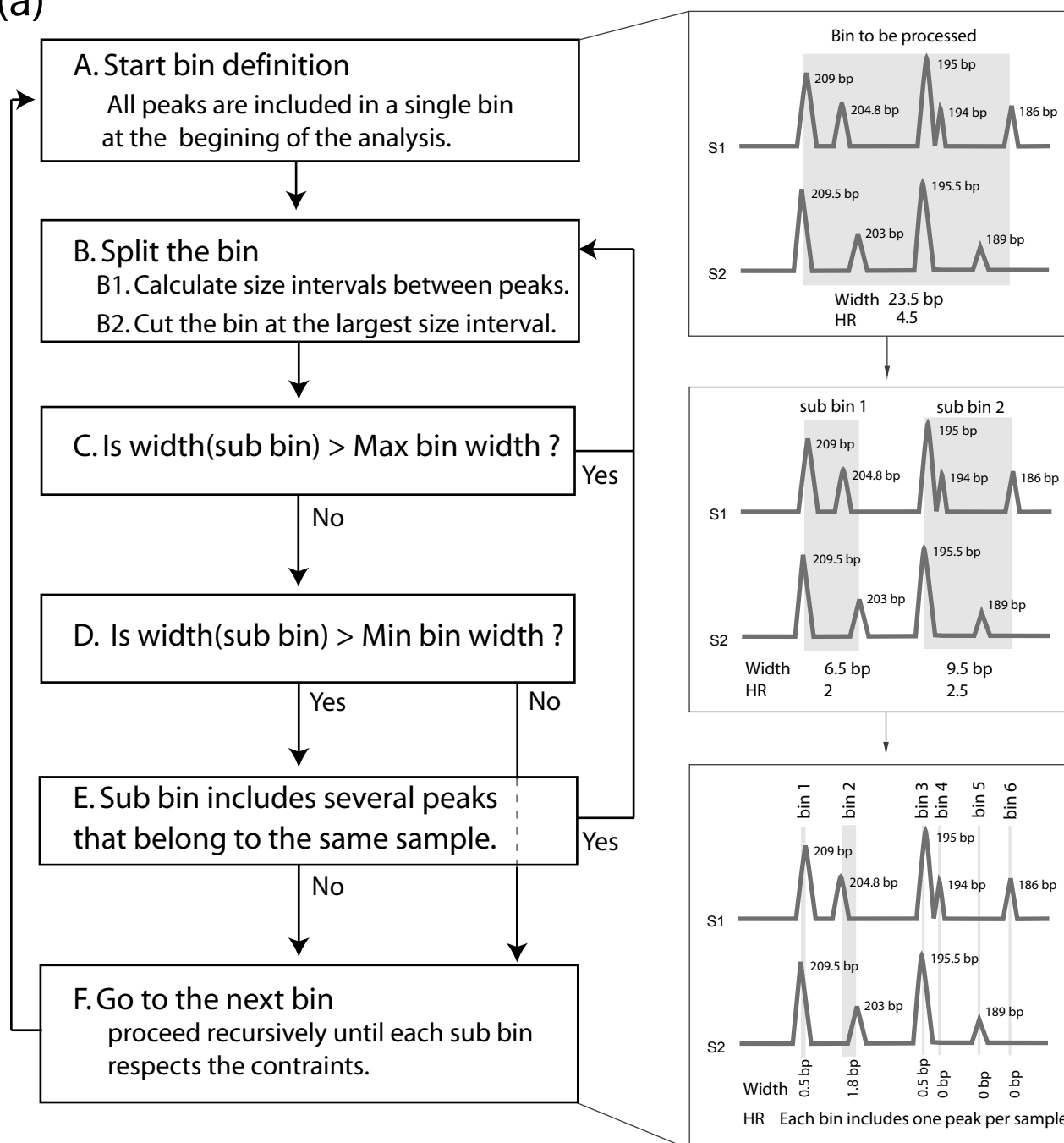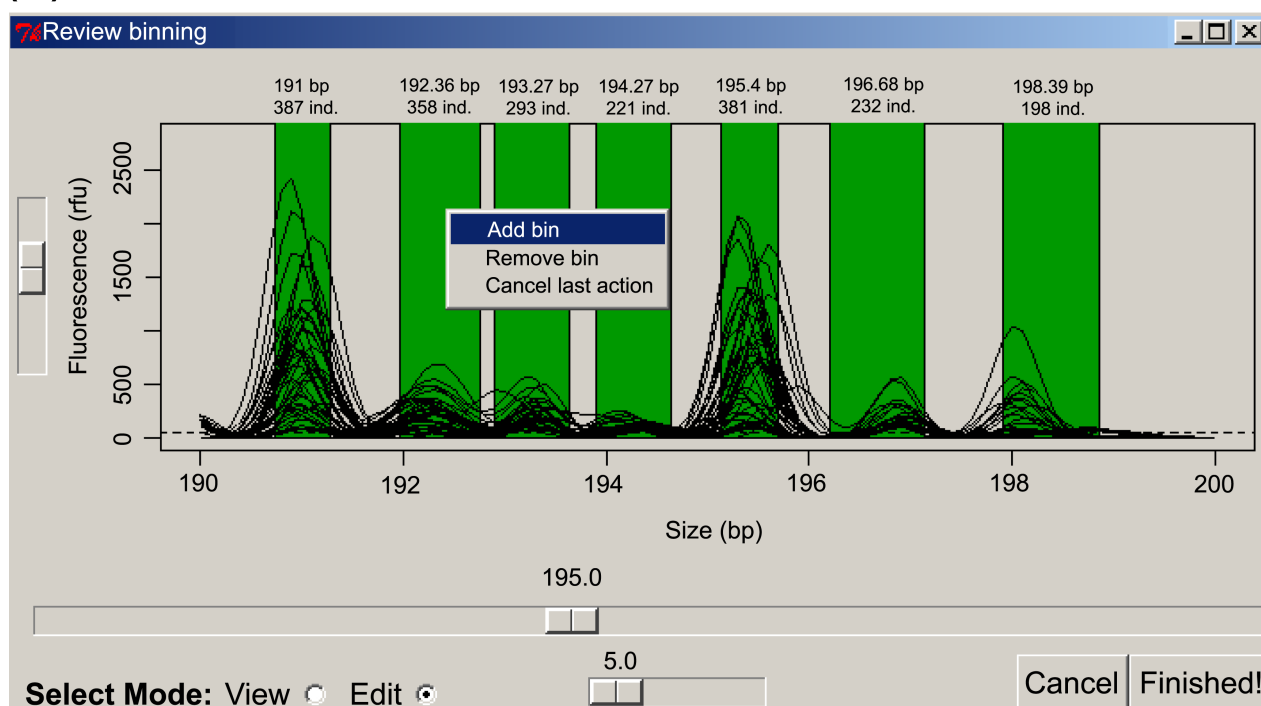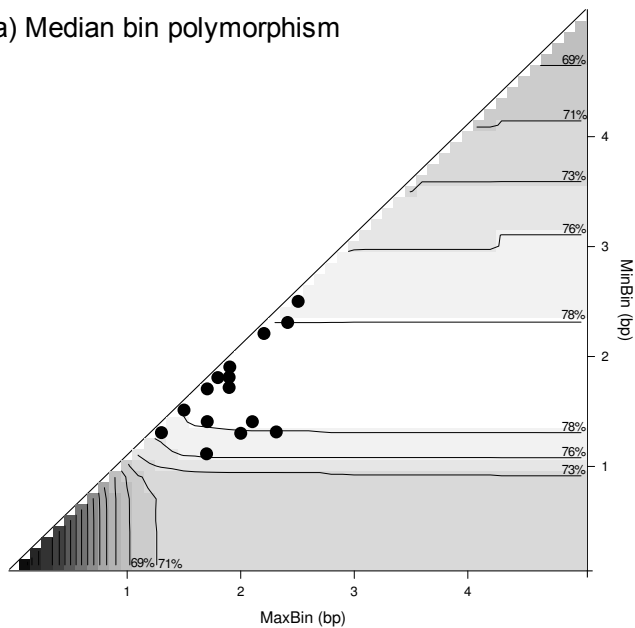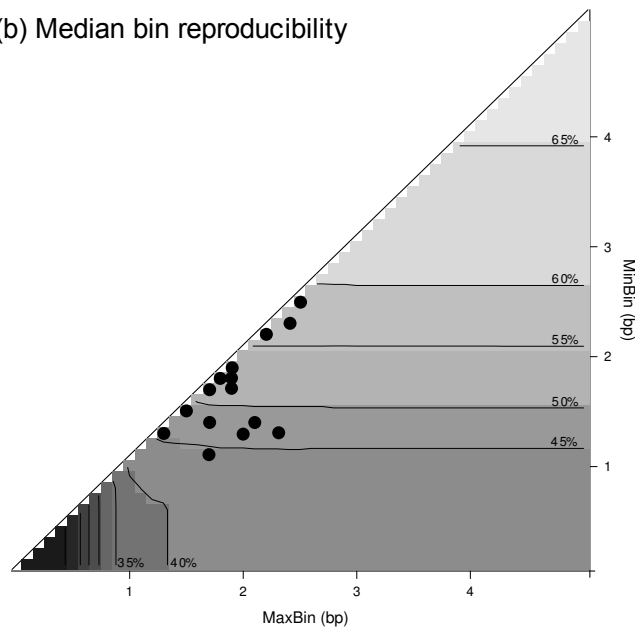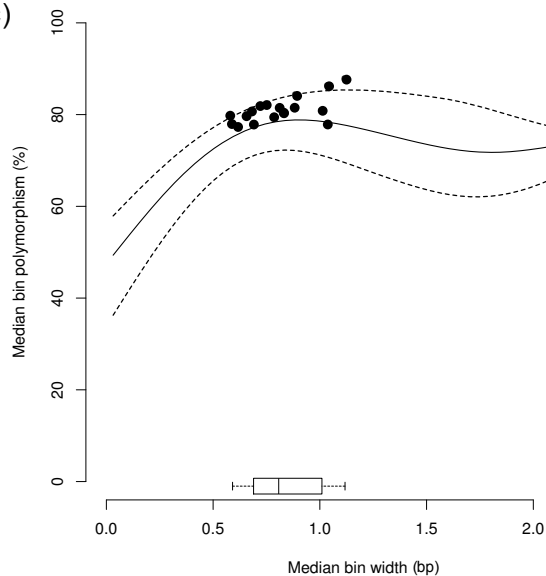