# Manual for EuroForMix v1

Author: Øyvind Bleka <Oyvind.Bleka.at.fhi.no>
Date: 12-04-2014

## (A) Installation and running program:

1) Run R (>=3.0.1) in Windows, Linux or MAC (http://cran.r-project.org/).
2) Required packages to run GUI:
   a. gWidgetstcltk (depends on digest,tcltk)
   b. gWidgets
3) Other required packages:
   a. cubature
      i. Required for multivariate integration (Integrated LR).
   b. forensim
      i. Required for qualitative Weight-of-Evidence.
4) Installation and run gammadnamix:
   a. install.packages("gammadnamix", repos="http://R-Forge.R-project.org")
   b. library(gammadnamix)
   c. euroformix()

## (B) GUI

## Sections:

0- Toolbar

1- Importing data

2- Model specification

3- MLE fit: ('Continuous LR (Maximum Likelihood based)')

4- Deconvolution (Deconvolution based on the continuous model)

5- Database Search (Database search based on the continuous and qualitative model )

6- Qual.LR (Qualitative model)

7- Generate data (Generation from the continuous model)

# 0. <u>Toolbar</u>

- File

    - o **Set directory**: The user may select the working directory of the R-program.

    - o **Open project**: The user may open an earlier project which is saved in a file on the form "projectname.Rdata".

    - o **Save project**: The user may save the existing project into a file with name "projectname".
        - Extension .Rdata is added automatically to project name.
        - All data imported to the program and resulting calculations are stored into a single project-file which may be open at any time in the program.
        - Saving a project makes:
            - Big reference databases are stored efficiently (the required space for the database is drastically reduced).
            - Time-consuming calculations are restored instantly (only required to be calculated ones).

    - o **Quit project**: When pushed, the user get question about saving project before terminating the GUI.

- Frequencies

    - o **Set size of frequency database**: User may specify number of samples 'N' used to create the population frequencies.
        - When new alleles from imported files are found, these are assigned as freq0.
            - If N=0 (this is default), freq0 is equal minimum observed frequency.
            - If N>0, freq0='5/(2N)'.
        - New alleles are updated to the population frequencies when:
            - When a reference database is imported.
            - When interpretations are done.
                - o Deconvolution, Weight-of-Evidence and 'Database search'
            - Frequencies are normalized for each of these two cases.
                - o **WARNING**: Normalizing may be done twice if new alleles (not seen in population frequency table or reference database) are observed in the evidence/reference profile.

    - o **Set number of wildcards in false positive match**: The user may specify number of wildcards in the random match probability statistics, which are applied when the user has imported and selected an evidence stain together with the population frequencies.

- Optimization

  o **Set number of random startpoints**: The user may set required number of independent random startpoints in the optimizer to ensure that the global maximum is attained for the Maximum Likelihood Estimator (MLE). Default is 3.

  o **Set variance of randomizer**: The user may set the variance parameter used for the random generation of startpoints used in optimizer. Default is 10**.**

- MCMC (Markov Chain Monte Carlo)

  o **Set number of samples**: The user may set the number of samples drawn from the posterior distribution of the parameters. Default is 10000.

  o **Set variance of randomizer**: The user may set the variance parameter scalar used in the 'Markov Chain Monte Carlo (MCMC) random walk Metropolis'. See vignette for details. Default is 10.
    - Note that this value should be tweaked such that acceptance rate of sampler are around 0.2 (to ensure global exploration in the parameter space).

- Integration

  o **Set relative error requirement**: The user may set the required estimated relative error used in the integration function adaptIntegrate {cubature}. See vignette for details. Default is 0.005.

  o **Set maximum of mu-parameter**: The user may set upper limit of mu-parameter (amount of DNA ). See vignette for details. Default is 20000.

  o **Set maximum of sigma-parameter**: The user may set upper limit of sigma-parameter (coefficient of variation). See vignette for details. Default is 1.

  o **Set maximum of stutter ratio-parameter**: The user may set upper limit of the stutter ratio parameter (xi). Default is 1.

- Deconvolution

  o **Set required summed probability**: The user may set required summed posterior genotype-probability which the deconvoluted list is ensured to contain. Default is 0.9999.

  o **Set max listsize:** The user may set maximum number of elements in the deconvoluted list. Default is 1000.
    - The greater max listsize, the more time-consuming (and memory consuming) the search-algorithm behind will be.

130    - Database search

131

132        o **Set maximum view-elements**: The user may set maximum number of individuals to
133          show from the reference-database. Default is 10000.
134            ▪ The greater 'value', the more time-consuming will it become to show table on
135              screen.
136            ▪ Note that the result table from the database search shows only the top 'value'-
137              ranked elements.

138

139        o **Set drop-in probability for qualitative model**: When searching database with
140          continuous LR model, the qualitative LR model is also considered with a specific drop-
141          in probability parameter given here (default is 0.05).

142

143    - Qual LR

144

145        o **Set upper range for sensitivity**: The user may specify the maximum allele dropout-
146          probability in the sensitivity plot (for a qualitative model). Default is 0.6.

147

148        o **Set nticks for sensitivity**: The user may specify number of grids of the allele dropout-
149          probability in the sensitivity plot (for a qualitative model). Default is 31.

150

151        o **Set required samples in dropout distr.**: The user may specify number of required
152          allele drop-out probability samples used to estimate the quantiles or meadian for the
153          distribution of the '*allele drop-out probability given number of observed alleles*'.

154

155        o **Set significance level in dropout distr.**: The user may specify the significance level in
156          the conservative LR calculation (i.e. the quantile for the distribution of the '*allele drop-
157          out probability given number of observed alleles*'). Default is 0.05.

158

159        o **Set number of tippets**: The user may specify number of random man tippet samples.
160          Default is 1e6.

161

162

163

164

165

166

167

168

169

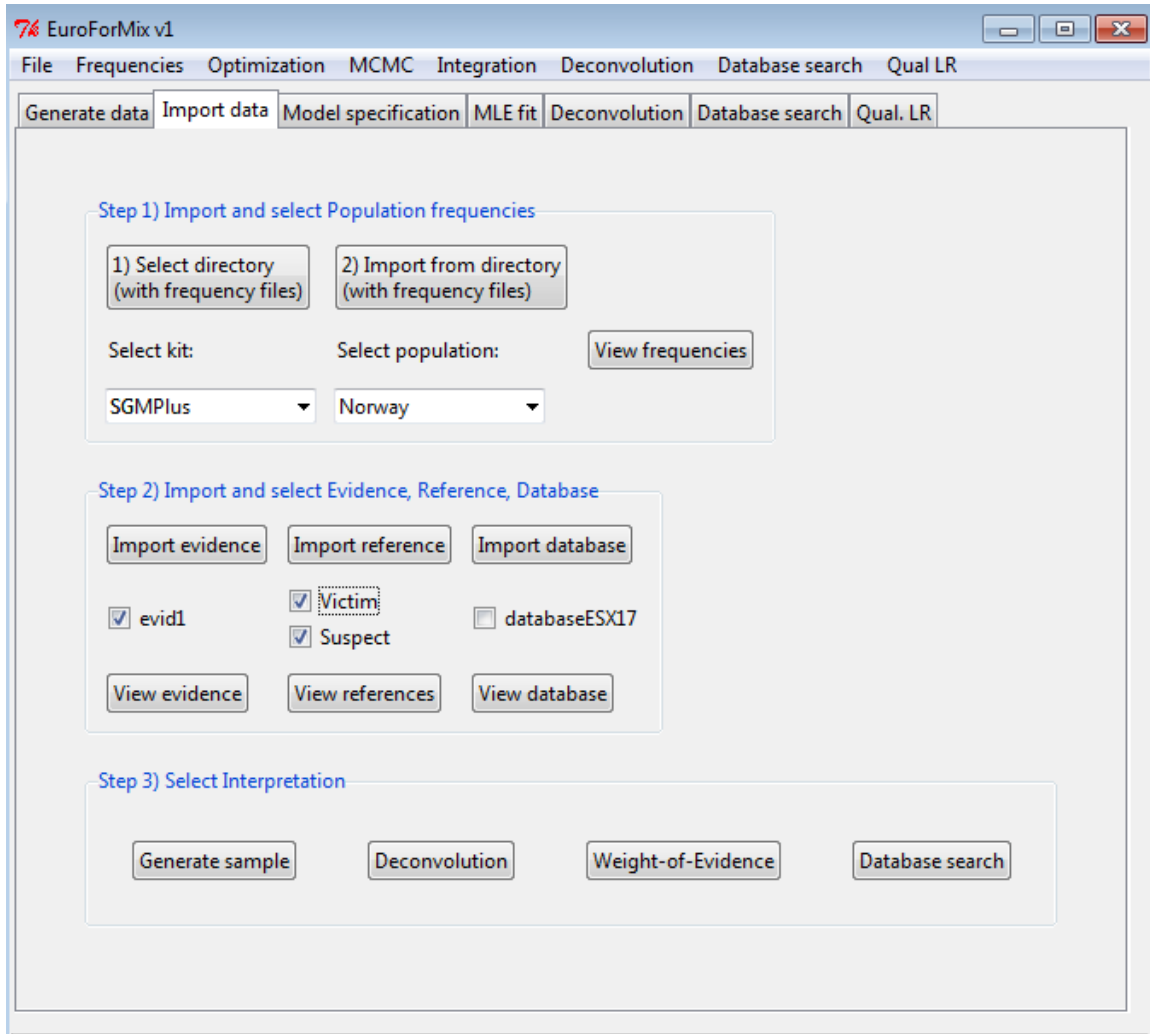170

171

# 1. <u>Importing data</u>

172

173



Figure 1: The figure shows the <u>Import data</u> GUI page where the user can import population frequencies, evidence stains, reference profiles and reference databases.

174
175
176
177
178

DATA IMPORT:

179
180

- **Common** for all files:

181
182
        o The extension (denotes file-type) of the file names does not matter. It may also have no extension at all.
183
184        o All imported files must be either comma, semi-colon or tab-separated (',',';','\t').
185        o Required/optional headers (all are capital invariant):
186            ▪ "**sample**" is required header for sample(s) name(s).
187                • The sample names are NOT capital invariant.
188                • If more than one header name contains "**sample**", it will select the header
189                  name which in addition contains "**name**" in the same string.

190 ▪ "**marker**" is required header for marker name(s).
191      • Marker names are capital invariant.
192      • If no header is found, the header containing "**loc**" will be used if found.
193 ▪ "**allele**" is required header(s) for allele-information.
194      • This may be a vector ("alleleX1",…,"allelleX10") of any length denoting
195      allele(s) to a given marker for a given sample. Here X1,…,X10 can be
196      anything.
197 ▪ "**height**" optional header(s) for peak height-information.
198      • This may be a vector ("heightX1",…,"heightX10") of any length
199      denoting peak height to the corresponding allele(s) in "allele". Here
200      X1,…,X10 can be anything.
201 ○ Note:
202 ▪ The imported data will use upper-letter of marker-names found in the file.
203 ▪ All imports are printed out in the terminal (see figure 2). From this, the user may
204 check that the data are imported correctly.
205

```
[1] "Raw fil import:"
   Sample.Name  Marker Allele.1 Allele.2 Allele.3 Allele.4 Allele.5 Allele.6 Height.1
1        evid1    AMEL        X        Y       NA       NA       NA       NA     2136
2        evid1 D3S1358       14       15     16.0       NA       NA       NA      178
3        evid1    TH01        6        7      9.3       NA       NA       NA      419
4        evid1  D21S11       27       29       NA       NA       NA       NA     1128
5        evid1  D18S51       15       17       NA       NA       NA       NA      467
6        evid1 D2S1338       17       19     20.0       23       NA       NA      290
7        evid1 D16S539        9       10     11.0       12       NA       NA      217
8        evid1     vWA       14       15     17.0       NA       NA       NA     1250
9        evid1 D8S1179       10       13     14.0       15       NA       NA      206
10       evid1     FGA       21       22       NA       NA       NA       NA      664
11       evid1 D19S433       13       14     15.2       NA       NA       NA     1157
   Height.2 Height.3 Height.4 Height.5 Height.6   ADO UD1  X
1      1015       NA       NA       NA       NA false  NA NA
2      2405     1982       NA       NA       NA false  NA NA
3       282     1871       NA       NA       NA false  NA NA
4      1750       NA       NA       NA       NA false  NA NA
5       524       NA       NA       NA       NA false  NA NA
6       619      259      649       NA       NA false  NA NA
7       312      743      619       NA       NA false  NA NA
8       440     1232       NA       NA       NA false  NA NA
9       352      978      827       NA       NA false  NA NA
10      714       NA       NA       NA       NA false  NA NA
11      781      922       NA       NA       NA false  NA NA
```

206
207 Figure 2: The figure shows the table format in the importing evidence stain file.
208
209 - **Import population frequencies**:
210
211 ○ Requires an own folder (population-folder) with **only** frequency-files.
212 ○ File-format:
213 ▪ Filename:
214      • The name of the filenames **needs** to be on the form "kit_population.ext",
215      where ext can be any extensions (or be missing as well).
216      • kit="kit-name" and population="population name"
217      • The kit-name must be consistent with the short-name of the kit
218      instrument. See ?plotEPG for more details.
219 ▪ File:
220      • First column needs to be allele-information (header-name may be
221      anything).

222          •  Other columns are frequency-information (header-name denotes the
223             locus name (loci names are converted to capital letters)).
224        o  To import frequencies:
225          ▪  Push "**1) Select directory**" button to select the population-folder with the
226           population frequency files.
227          ▪  Push "**2) Import from directory**" button to import the population frequency
228           files from the selected folder.
229             •  It is possible to **add new files** into the selected population-folder **at any**
230              **time** and push the button once again to include new information to the
231              dropdown-list.
232        o  Selection of kit and population:
233          ▪  After importing the frequency-files (after pushed (**2**)), the user may select
234           wanted kit and population from the two drop down lists at any time* (*not after
235           a reference-database file has been imported).
236             o  This can be useful to see the EPG layout for different selected
237              kits.
238
239   -  **Import Evidence/Reference** sample (see figure 2 and figure 3):
240
241        o  **Multiple** evidence or reference profiles are **allowed** in each file.
242        o  In evidence files:
243          ▪  "height" header is required for analysis Deconvolution, Weight-of-Evidence
244           (continuous model) and 'Database search'. For 'Qualitative LR' this is not
245           required.
246        o  In reference files:
247          ▪  "height" header is optional but will not be used further in any analysis.
248        o  Note:
249          ▪  The import function will not check:
250             •  That the length of allele and heights are equal long for a given locus.
251          ▪  Loci without any allele-information (i.e. empty or dropped out), are **NOT**
252           imported.
253

```
[1] "Raw fil import:"
  SampleName  Marker Allele1 Allele2
1       Victim D3S1358   16.0   15.0
2       Victim    TH01    9.3    9.3
3       Victim  D21S11   29.0   27.0
4       Victim  D18S51   17.0   15.0
5       Victim D2S1338   23.0   19.0
6       Victim D16S539   11.0   12.0
7       Victim     VWA   14.0   17.0
8       Victim D8S1179   14.0   15.0
9       Victim     FGA   22.0   21.0
10      Victim D19S433   13.0   15.2
11     Suspect D3S1358   16.0   15.0
12     Suspect    TH01    6.0    7.0
13     Suspect  D21S11   29.0   35.0
14     Suspect  D18S51   11.0   14.0
15     Suspect D2S1338   17.0   20.0
16     Suspect D16S539    9.0   10.0
17     Suspect     VWA   15.0   17.0
18     Suspect D8S1179   10.0   13.0
19     Suspect     FGA   22.0   25.0
20     Suspect D19S433   14.0   14.0
```

254
255          Figure 3: The figure shows the table format in the importing reference file.

256
257 - **Import Reference Database (**see figure 4):
258      o Exactly same format as reference files.
259      o Multiple database file may be imported (**must** be done one-at-the-time).
260      o **Requires** that population frequencies are imported and selected.
261          ▪ **WARNING**: Population frequencies may not be changed again after database
262          importing!
263      o Note:
264          ▪ The ranking of databases are done over all selected databases.
265          ▪ Same samples within a database needs to be in same block but markers within
266          sample can be different orders.
267          ▪ Some samples **may** have more/less markers than others (e.g. SGMplus profiles
268          contra ESX17).
269             • **Missing markers** for a sample are given with NA.
270          ▪ Only markers shared with selected population frequencies are imported.
271             • The imported database files may contain different markers.
272          ▪ Homozygote genotype may have an empty allele under 'Allele 2'.
273          ▪ The database file may contain **any** number of individuals.
274      o Tips:
275          ▪ It is more efficient to import several small databases than one big.
276             • Time usage to import a database file with 16 markes:
277                 o 1e6 profiles takes about 130 seconds
278                    ▪ Requires ~1.3GB memory
279                 o 5e6 profiles takes about 800 seconds.
280                    ▪ Requires ~6.1GB memory
281          ▪ Save a lot of time and memory by storing a project to file (See File under
282          toolbar). The imported database will be stored very efficiently.
283

```
[1] "Raw fil import:"
                          Sample.Name    Marker Allele.1 Allele.2
1   00-JP0001-14_20142342311_NO-3241   D3S1358       14       15
2   00-JP0001-14_20142342311_NO-3241      TH01        7      9.3
3   00-JP0001-14_20142342311_NO-3241    D21S11       29       30
4   00-JP0001-14_20142342311_NO-3241    D18S51       13       17
5   00-JP0001-14_20142342311_NO-3241  D10S1248       12       13
6   00-JP0001-14_20142342311_NO-3241   D1S1656       11       14
7   00-JP0001-14_20142342311_NO-3241   D2S1338       17       19
8   00-JP0001-14_20142342311_NO-3241   D16S539       10       11
9   00-JP0001-14_20142342311_NO-3241  D22S1045       15       16
10  00-JP0001-14_20142342311_NO-3241       VWA       17       18
11  00-JP0001-14_20142342311_NO-3241   D8S1179       12       13
12  00-JP0001-14_20142342311_NO-3241       FGA       19       22
13  00-JP0001-14_20142342311_NO-3241    D2S441       11       10
14  00-JP0001-14_20142342311_NO-3241   D12S391       17       18
15  00-JP0001-14_20142342311_NO-3241   D19S433       13       14
16  00-JP0001-14_20142342311_NO-3241      SE33       15       21
17  00-JP0001-14_20142342311_NO-3241      AMEL        X        Y
18  00-JP0002-14_20142342311_NO-3242   D3S1358       15       18
19  00-JP0002-14_20142342311_NO-3242      TH01        6        9
20  00-JP0002-14_20142342311_NO-3242    D21S11       28     31.2
21  00-JP0002-14_20142342311_NO-3242    D18S51       13       18
22  00-JP0002-14_20142342311_NO-3242  D10S1248       13       13
23  00-JP0002-14_20142342311_NO-3242   D1S1656       15     18.3
24  00-JP0002-14_20142342311_NO-3242   D2S1338       25       25
25  00-JP0002-14_20142342311_NO-3242   D16S539       11       13
26  00-JP0002-14_20142342311_NO-3242  D22S1045       15       16
27  00-JP0002-14_20142342311_NO-3242       VWA       14       17
```
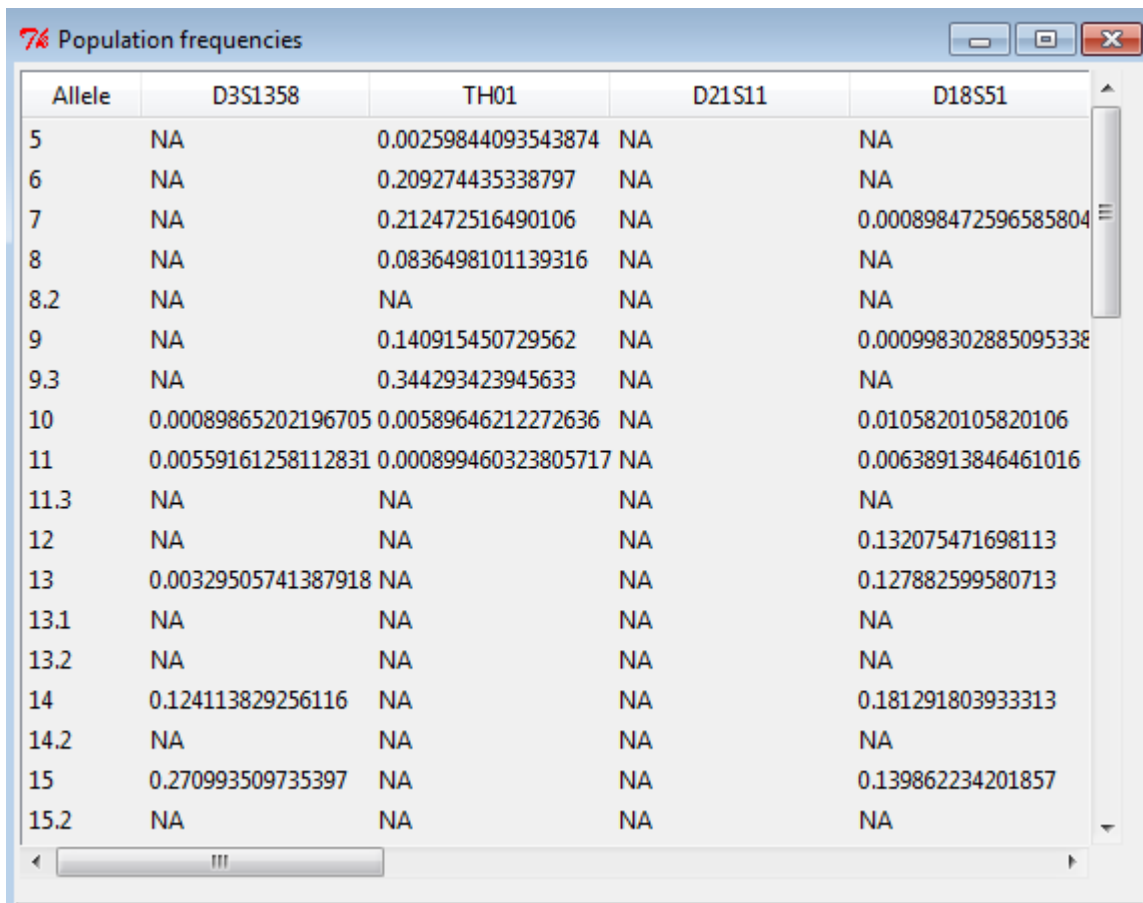284
285 Figure 4: The figure shows the table format in the importing reference database file.

VIEW DATA:
287
288     -    **View frequencies** (see figure 5 for the Norwegian SGMPlus population):
289
290          o   Creates a new window which shows the selected population frequencies in a table.
291          o   If any evidence profiles(s) are selected after evidence-import, the software makes a
292              'false positive probability' – plot for each selected profiles.
293              ▪   The plot (figure 6) shows the probability that a random individual (**'false
294                  positive probability'**) matching at least (2*n-wildcardsize) up to 2*n alleles
295                  (MAC) with a **selected evidence** profile. Here **n** is number of considered loci
296                  (which are both in evidence and population frequencies) and wildcardsize is
297                  number of allowed mismatches (default is wildcardsize =7).
298              ▪   wildcardsize can be changed under "Frequencies" in Toolbar by changing value
299                  **Set number of wildcards in false positive match.**
300          o   Note:
301              ▪   Only allele-information in evidence-profiles are used.
302              ▪   New alleles which are not found in the selected population are assumed to have
303                  allele-frequency 0.
304
305

| Allele | D3S1358 | TH01 | D21S11 | D18S51 |
|--------|---------|------|--------|--------|
| 5 | NA | 0.00259844093543874 | NA | NA |
| 6 | NA | 0.209274435338797 | NA | NA |
| 7 | NA | 0.212472516490106 | NA | 0.000898472596585804 |
| 8 | NA | 0.0836498101139316 | NA | NA |
| 8.2 | NA | NA | NA | NA |
| 9 | NA | 0.140915450729562 | NA | 0.000998302885095338 |
| 9.3 | NA | 0.344293423945633 | NA | NA |
| 10 | 0.00089865202196705 | 0.00589646212272636 | NA | 0.01058201058201106 |
| 11 | 0.00559161258112831 | 0.000899460323805717 | NA | 0.00638913846461016 |
| 11.3 | NA | NA | NA | NA |
| 12 | NA | NA | NA | 0.132075471698113 |
| 13 | 0.00329505741387918 | NA | NA | 0.127882599580713 |
| 13.1 | NA | NA | NA | NA |
| 13.2 | NA | NA | NA | NA |
| 14 | 0.124113829256116 | NA | NA | 0.181291803933313 |
| 14.2 | NA | NA | NA | NA |
| 15 | 0.270993509735397 | NA | NA | 0.139862234201857 |
| 15.2 | NA | NA | NA | NA |

306
307     Figure 5: The figure shows the viewed frequencies for the Norwegian SGMPlus frequencies.
308

Random man probability having number of allele matches>=k
Sample: evid1
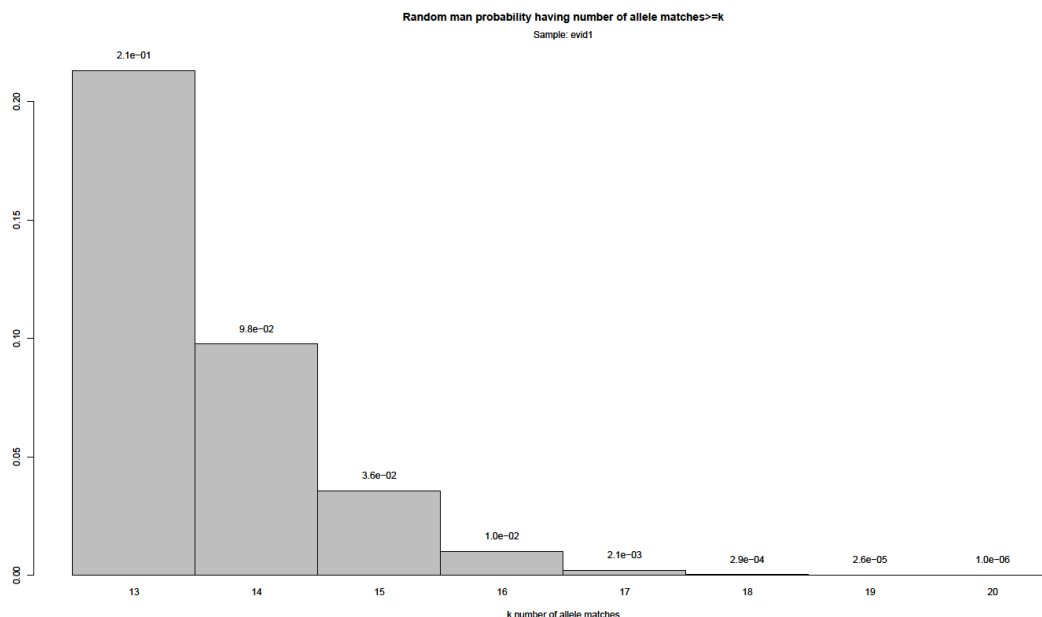
309
310  Figure 6: The figure shows the random match probability of matching with at least k number of alleles
311                    (in reference) with the observed alleles in evidence.
312
313      -   **View evidence** (for selected evidence):
314
315          o   Prints imported alleles (and peak heights if any) for each selected evidence profile(s)
316              (see figure 7).
317

```
[1] "Samplename: evid1"
           Allele           Height
AMEL      "X/Y"            "2136/1015"
D3S1358  "14/15/16"        "178/2405/1982"
TH01     "6/7/9.3"         "419/282/1871"
D21S11   "27/29"           "1128/1750"
D18S51   "15/17"           "467/524"
D2S1338  "17/19/20/23"     "290/619/259/649"
D16S539  "9/10/11/12"      "217/312/743/619"
VWA      "14/15/17"        "1250/440/1232"
D8S1179  "10/13/14/15"     "206/352/978/827"
FGA      "21/22"           "664/714"
D19S433  "13/14/15.2"      "1157/781/922"
```

318
319  Figure 7: The figure shows the printed alleles and heights in the imported evidence.
320
321          o   Plots EPG(s) (see figure 7) for each selected evidence profile(s)
322              ▪   Requires that user have imported "Population frequencies".
323              ▪   The kit selected under '**Select kit**' denotes the EPG format.
324              ▪   Loci in evidence which are **inconsistent** with the ones in selected kit (or
325                  missing) are **not shown** in plot.
326              ▪   Evidence profiles without peak heights for corresponding alleles are given with
327                  peak height equal 1.
328          o   Note:
329              ▪   See ?plotEPG to see which kit-formats that are supported.

330        ▪ Reference profiles can be imported as evidence profiles and shown in a EPG.
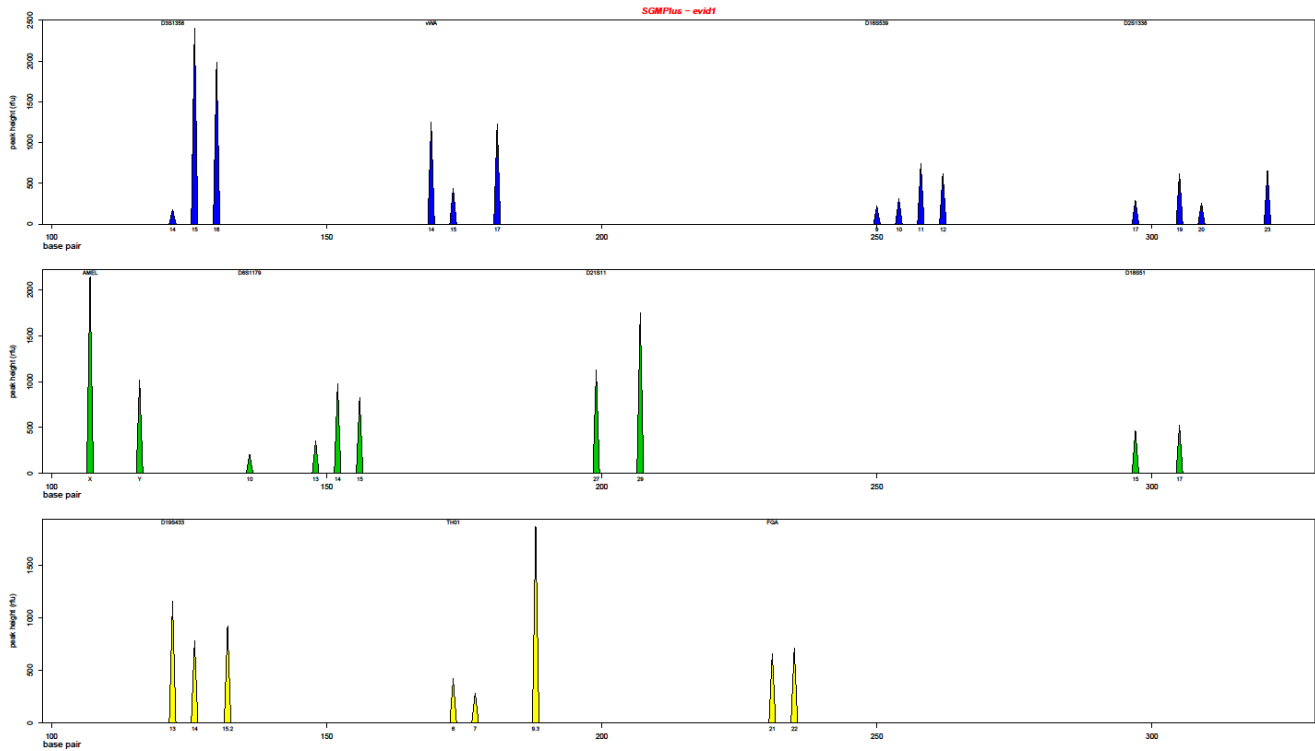331



332
333
334 Figure 8: The figure shows the plotted EPG (on selected SGMPlus kit format) of the imported evidence
335                                          stain.
336
337    -   **View reference** (for selected reference):
338        o Prints imported genotypes for each selected reference profile(s) (figure 9).
339        o If any evidence profiles(s) are selected after evidence-import, the software counts
340          number of matching alleles (MAC) for each loci of the selected reference profiles, for
341          each selected evidences (figure 10).
342            ▪ MAC = number of alleles for the reference which are included in the evidence.
343            ▪ nLocs = number of considered loci when counting MAC.
344

```
                         Victim          Suspect
          D3S1358   "16/15"         "16/15"
          TH01      "9.3/9.3"       "6/7"
          D21S11    "29/27"         "29/35"
          D18S51    "17/15"         "11/14"
          D10S1248  "15/13"         "13/13"
          D1S1656   "12/17.3"       "15/16"
          D2S1338   "23/19"         "17/20"
          D16S539   "11/12"         "9/10"
          D22S1045  "15/16"         "15/15"
          VWA       "14/17"         "15/17"
          D8S1179   "14/15"         "10/13"
          FGA       "22/21"         "22/25"
          D2S441    "10/14"         "11/11"
          D12S391   "18.3/22"       "18/19"
          D19S433   "13/15.2"       "14/14"
          SE33      "30.2/33.2"  "27.2/29.2"
```

345
346 Figure 9: The figure shows the printed alleles of the imported reference profiles.

```
[1] "Number of matching alleles with samplename evid1:
           Victim Suspect
AMEL           NA       NA
D3S1358         2        2
TH01            2        2
D21S11          2        1
D18S51          2        0
D2S1338         2        2
D16S539         2        2
VWA             2        2
D8S1179         2        2
FGA             2        1
D19S433         2        2
MAC            20       16
nLocs          10       10
```

Figure 10: The figure shows number of matching alleles and total (MAC) with the imported and selected evidence stain. By combining the observed MAC and figure 7, the random match probability of observing MAC is useful for providing an extended version of "Random man not excluded"-statistics: The random match probability for Victim (MAC=20) becomes 1/1000000, while only 1/100 for Suspect (MAC=16).

- **View database** (see figure 11 for selected database):

  o Creates a new window (for each selected database) which shows the genotypes for every reference in the database.
    ▪ "NA" means that the genotype of a reference was missing.
  o If any evidence profiles(s) are selected after evidence-import, the software counts number of matching alleles (MAC) for all references in the database against each of the selected evidences (see figure 12). The results are shown in a MAC-ranked table in a new window (for each selected database).
    ▪ **MAC** = total number of alleles for the reference which are included in the evidence.
      • Summed over all selected evidences.
    ▪ **nLocs** is number of reference-loci which has been used to evaluate the MAC.
  o Note:
    ▪ Max number of individuals to view in a database can be changed with selecting **Set maximum view-elements** under "Database search" in toolbar.

References in imported database databaseESX17

| Reference | D3S1358 | TH01 | D21S11 | D18S51 | D2S1338 | D16S539 | VWA | D8S1179 | FGA | D19S433 |
|---|---|---|---|---|---|---|---|---|---|---|
| 00-JP0001-14_20142342311_NO-3241 | 14/15 | 7/9.3 | 29/30 | 13/17 | 17/19 | 10/11 | 17/18 | 12/13 | 19/22 | 13/14 |
| 00-JP0002-14_20142342311_NO-3242 | 15/18 | 6/9 | 28/31.2 | 13/18 | 25/25 | 11/13 | 14/17 | 12/12 | 21/23 | 12/14.2 |
| 00-JP0003-14_20142342311_NO-3243 | 16/18 | 9.3/9.3 | 30/30 | 13/18 | 17/18 | 8/12 | 16/18 | 12/13 | 18/24 | 14/15 |
| 00-JP0004-14_20142342311_NO-3244 | 18/18 | 7/9.3 | 29/32.2 | 12/22 | 19/23 | 11/11 | 14/16 | 13/13 | 20/20 | 13.2/14 |
| 00-JP0005-14_20142342311_NO-3245 | 15/17 | 7/8 | 28/33.2 | 12/17 | 19/25 | 13/13 | 17/18 | 12/14 | 20/21 | 15/15 |
| 00-JP0006-14_20142342311_NO-3246 | 14/18 | 7/9.3 | 28/32.2 | 11/15 | 20/24 | 9/13 | 15/16 | 13/13 | 22/22 | 14/15 |
| 00-JP0007-14_20142342311_NO-3247 | 15/19 | 9.3/9.3 | 30/32 | 14/19 | 17/23 | 9/10 | 16/16 | 10/12 | 23/25 | 13/15 |
| 00-JP0008-14_20142342311_NO-3248 | 14/16 | 9/9.3 | 30/30.2 | 14/18 | 17/23 | 9/11 | 16/18 | 13/14 | 20/20 | 12/14 |
| 00-JP0009-14_20142342311_NO-3249 | 14/16 | 7/7 | 30/30 | 12/16 | 21/22 | 12/12 | 14/16 | 12/13 | 21/21 | 12/12 |
| 00-JP00010-14_20142342311_NO-3241 | 15/16 | 6/6 | 30/32 | 16/17 | 21/23 | 9/14 | 18/18 | 13/15 | 19/22 | 13/14 |
| 00-JP00011-14_20142342311_NO-3241 | 15/17 | 6/9 | 29/30 | 15/16 | 17/25 | 12/12 | 15/20 | 12/13 | 21/23 | 15/15 |
| 00-JP00012-14_20142342311_NO-3241 | 15/17 | 7/9.3 | 30/31.2 | 14/19 | 19/20 | 10/12 | 17/17 | 13/15 | 20/24 | 12/14 |
| 00-JP00013-14_20142342311_NO-3241 | 17/18 | 6/9 | 28/29 | 12/19 | 17/24 | 11/13 | 17/17 | 11/13 | 22/25 | 14/14 |
| 00-JP00014-14_20142342311_NO-3241 | 15/18 | 9/9.3 | 29/30 | 13/18 | 18/24 | 9/13 | 16/16 | 12/14 | 21/24 | 15/15 |
| 00-JP00015-14_20142342311_NO-3241 | 16/16 | 8/9.3 | 30/30 | 12/15 | 17/24 | 9/11 | 15/16 | 11/14 | 19/23 | 13/15 |
| 00-JP00016-14_20142342311_NO-3241 | 14/15 | 6/9.3 | 28/31 | 15/17 | 23/25 | 11/12 | 14/14 | 12/13 | 20/21 | 13/14 |
| 00-JP00017-14_20142342311_NO-3241 | 17/18 | 6/7 | 29/33.2 | 13/14 | 19/19 | 13/13 | 14/16 | 12/13 | 18/24 | 14/15 |
| 00-JP00018-14_20142342311_NO-3241 | 15/20 | 6/7 | 29/30 | 15/7 | 17/17 | 9/13 | 14/17 | 12/14 | 20/26 | 13/15 |
| 00-JP00019-14_20142342311_NO-3241 | 15/18 | 7/7 | 28/29 | 13/16 | 17/25 | 12/12 | 17/17 | 11/14 | 20/21 | 14/14 |
| 00-JP00020-14_20142342311_NO-3242 | 16/16 | 7/9.3 | 29/29 | 16/19 | 17/24 | 11/11 | 16/17 | 11/13 | 19/23 | 13/14 |
| 00-JP00021-14_20142342311_NO-3242 | 14/14 | 9/9 | 29/30 | 13/19 | 22/24 | 9/12 | 14/18 | 13/14 | 19/20 | 14/16 |
| 00-JP00022-14_20142342311_NO-3242 | 15/17 | 6/8 | 29/31.2 | 14/18 | 17/18 | 11/11 | 18/18 | 13/13 | 20/20 | 15/15 |
| 00-JP00023-14_20142342311_NO-3242 | 14/16 | 7/7 | 31.2/32.2 | 13/14 | 20/23 | 11/11 | 14/16 | 13/14 | 21/19.2 | 14/15 |
| 00-JP00024-14_20142342311_NO-3242 | 15/17 | 7/9.3 | 30/31.2 | 15/17 | 20/24 | 11/12 | 16/16 | 15/15 | 21/21 | 14/14.2 |
| 00-JP00025-14_20142342311_NO-3242 | 16/17 | 6/7 | 28/29 | 14/16 | 17/19 | 11/12 | 14/17 | 14/14 | 22/24 | 13/14 |

Figure 11: The figure shows the viewed references inside the imported ESX17 database which are presented only with SGMPlus profiles since the selected kit for the imported frequencies was SGMPlus_Norway.

Number of sample matching alleles in refe...

| Reference | evid1 | nLocs |
|---|---|---|
| 00-JP00059-14_20142342311_NO-32459 | 17 | 10 |
| 00-JP0001-14_20142342311_NO-3241 | 15 | 10 |
| 00-JP00016-14_20142342311_NO-32416 | 15 | 10 |
| 00-JP00025-14_20142342311_NO-32425 | 15 | 10 |
| 00-JP00066-14_20142342311_NO-32466 | 15 | 10 |
| 00-JP00036-14_20142342311_NO-32436 | 14 | 10 |
| 00-JP00057-14_20142342311_NO-32457 | 14 | 10 |
| 00-JP00019-14_20142342311_NO-32419 | 13 | 10 |
| 00-JP00020-14_20142342311_NO-32420 | 13 | 10 |
| 00-JP00023-14_20142342311_NO-32423 | 13 | 10 |
| 00-JP00024-14_20142342311_NO-32424 | 13 | 10 |
| 00-JP00033-14_20142342311_NO-32433 | 13 | 10 |
| 00-JP00042-14_20142342311_NO-32442 | 13 | 10 |
| 00-JP00049-14_20142342311_NO-32449 | 13 | 10 |

Figure 12: The figure shows the sorted references (in the reference database) with respect to MAC (total number of matching alleles) to the selected evidence.

<div align="center">INTERPRETATIONS:</div>

- **Generate sample**:

  - Generates alleles using the population frequencies and draws peak heights for a specified hypothesis using the continuous model as described in the vignette.
  - Requires: Imported population frequencies.
  - Feature: Allele drop-out, Drop-in (with a peak height model) and stutter.

- **Deconvolution**:

  - Deconvolution ranks the most probable combined genotype profiles given **a specified hypothesis** and the Maximum Likelihood Estimates of the parameters in the continuous model (as given in the vignette).
  - Requires: Imported population frequencies and selection of at least one evidence profile with peak height information. References are optional to condition on in the hypothesis.
  - Feature: Model may handle replicates, allele drop-in, drop-out and stutter.

- **Weight-of-Evidence**:

  - Weight-of-Evidence is done by comparing the Likelihood Ratio (LR) between the specified hypotheses Hp (prosecution) and Hd (defence) using the continuous model as given in the vignette.
  - Modules:
    - 1) 'Continuous LR' (Maximum Likelihood based)
      - Optimizes (maximum) the model parameters in the continuous model.
    - 2) 'Continuous LR' (Integrated Likelihood based)
      - Integrates out the model-parameters in the continuous model.
    - 3) 'Qualitative LR' (semi-continous)
      - Explores LR as a function of allele dropout probability parameter.

  - Requires:
    - Imported population frequencies, **at least one** evidence profile and **at least one** reference profile (suspect) to weight evidence for. Additional reference profiles are optional to condition on in the hypotheses.
    - 'Continuous LR' requires evidence(s) including peak heights, 'Qualitative LR' only requires allele data.
  - Feature:
    - The continuous model: Handles replicates, allele drop-in, drop-out, stutter and fst-correction.
    - The semi-continuous model: Handles replicates, allele drop-in, drop-out and fst-correction.

428    -    **Database search:**

429

430         o   Does weight-of-evidence by comparing the Likelihood Ratio (LR) between the specified
431             hypotheses Hj (reference j in database) and Hd (defence) using the continuous model as
432             given in the vignette.
433         o   Modules:
434             ▪   1) 'Continuous LR' (Maximum Likelihood based)
435             ▪   2) 'Continuous LR' (Integrated Likelihood based)
436             ▪   3) 'Qualitatitve LR' (Semi-continuous based)
437         o   Requires: Imported population frequencies, **at least one** evidence profile with **peak
438             height** information and **at least one** reference-database. Reference profiles are optional
439             to condition on in the hypotheses.
440         o   Feature: Model may handle replicates, allele drop-in, drop-out, stutter and fst-correction.
441         o   The continuous LR value is showed together with qualitative LR and MAC.

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

# 2. Model specification

Figure 13: The figure shows the <u>Model Specification</u> GUI page for **Weight-of-Evidence** based on Likelihood Ratio calculation.

- **Evidence(s)**:

    o Shows selected evidence(s) from 'Import data'.
    o All interpretations support **multiple replicates**.
        ▪ Note: All replicates are assumed to have same parameter sets.

- **Contributors under Hp**

    o Case: **Weight-of-Evidence** or **'Database search'):**
        ▪ User may condition on selected references (from 'Import data') in the hypothesis Hp.
        ▪ #unknowns under Hp: Denotes number of unknown contributors under the prosecution hypothesis Hp.
    o Case: **'Database search':**
        ▪ The individual in the reference-database is already included in the hypothesis Hp.
    o Case: **Deconvolution** or **'Generate sample'**:
        ▪ This block is not considered, since Deconvolution only considers the model under Hd, and sample generation is done only under a specific hypothesis.

- **Contributors under Hd** (same for **all** cases):

    o User may condition on selected references (from 'Import data') in the hypothesis Hd.
    o #unknowns under Hd: Denotes number of unknown contributors under the prosecution hypothesis Hd.
    o Case: **Weight-of-Evidence** or **'Database search'**:
        ▪ References which are conditioned under Hp but not under Hd, will be assumed to be a '**known non-contributor'** under Hd (this is relevant when fst>0).

- **Continuous Model Parameters** and **Qualitative Model Parameters**:

    o The Continuous Model Parameter section is only used for "Continuous LR" Calculations, while Qualitative Model Parameters section is only used for 'Qualitative LR' Calculations.

    o '**Probability of drop-in**': [0,1]
        ▪ Assumed probability of a random allele drop-in to the evidence at a given locus. See vignette for more details.
        ▪ This is default 0 for continuous models and 0.05 for qualitative models.

    o **fst-correction**: [0,1]
        ▪ Assumed co-ancestry parameter assigned in the genotype probability for each contributor in the hypotheses. See vignette for more details.
        ▪ This is default 0 for continuous models and 0.02 for qualitative models.

521    o   Case **'Database search':**
522        ▪   When doing database search with "Continuous LR" Calculations, the allele drop-
523            in probability for the qualitative LR can be changed by **Set drop-in probability**
524            **for qualitative model** under "Database search" in toolbar (default is 0.05).
525            When doing database search with "Qualitative LR" Calculations, this value is
526            ignored in favor of the specification under "Qualitative Model Parameters".
527
528    o   Case **Generation** and **Deconvolution**:
529        ▪   The Qualitative Model Parameters section is removed.
530
531   -   **Advanced Parameters**
532
533        o   **Q-assignation**:
534
535            ▪   If checked, all alleles **not** present in the evidence are considered as allele "99".
536                Its frequency will be given as the sum of the frequencies for all the "non-
537                present" alleles.
538            ▪   If unchecked, the original alleles in the population are used as before.
539
540        o   '**Detection threshold'**: [0,->)
541
542            ▪   The threshold of required allele peak heights of whether an allele is present in
543                the evidence or not.
544                •   Note: If peak heights in evidence are lower than the specified threshold,
545                    the corresponding alleles (and peak heights) below threshold **are**
546                    **removed** automatically. This may cause some loci to become empty.
547
548        o   '**Stutter ratio'**: [0,1]
549
550            ▪   Only used for 'Continuous LR' Calculations.
551            ▪   Stutter ratio is a constant parameter "**xi**" which denotes the proportion of peak
552                heights from allele 'a' which is added to allele 'a-1'. See vignette for more
553                details.
554                •   If allele 22 with peak height y_22 is contributed by a contributor and
555                    allele 23 did not have any observed peak height, then the stutter
556                    contribution to allele 21 from allele 22 will be (**xi** * y_22).
557
558        o   '**Dropin peak height hyperparam'**: [0,1]
559
560            ▪   Only used for 'Continuous LR'.
561            ▪   Assumed hyper-parameter to model the peak height of the dropped in allele
562                caused by a 'random allele drop-in' if '**Probability of drop-in'**>0. See vignette
563                for more details.
564
565   -   **'Database(s) to search'** (case**: 'Database search'**)
566        o   Lists the selected imported reference-database(s) to do the database search for.

<div align="center">DATA SELECTION</div>

- **Select/unselect loci**:

    o  The user may select or unselect loci for each selected evidence(s) and reference(s) from "Import data"
    o  If a locus has been unselected for any of the evidence(s) or reference(s), the unselected locus will not be evaluated at all.
    o  Note: Evidence with more than 30 loci will not be able to be selected.

- **Missing data**:

    o  Data with missing allele in any of the loci will automatically be deselected (inactivated) such that the corresponding loci will be unavailable to evaluate.
    o  For continuous LR evaluation:
        ▪ If peak heights (in any of the evidence(s)) are missing for any selected locus, the user gets a message about deselecting the issued loci before proceeding.

- **New alleles**:

    o  If new alleles (does not exist in the population frequency table) occurs in the imported evidence or reference profile, the new alleles are assigned with allele frequency 'freq0'. 'freq0' is equal minimum observed frequency in population if N=0, or 'freq0'=5/(2N) where N is size of imported frequency database under "Frequencies" in Toolbar. The frequencies are after normalized.


<div align="center">SHOW SELECTED DATA</div>

- **Plot EPG:**

    o  **Prints** the selected evidence sample(s), reference(s) and considered population frequencies which are eventually used for further analysis **out to terminal**.
    o  The selected evidence samples are shown in an EPG-plot.
        ▪ Note: Alleles with corresponding peak heights below the specified "Detection Threshold" are removed.


<div align="center">CALCULATIONS</div>

- **'Continuous LR (Maximum Likelihood based) '** (case **Weight-of-Evidence** and **'Database search')**:

    o  Maximizes the Likelihood of the unknown parameters in the continuous model given the assumed model so they attain maximum values for the specified hypothesis Hd (and Hp in case of Weight-of-Evidence).

613 ▪ The optimizer should return a global maximum. However, it may sometimes just
614 return a local maximum. Number of start-points should be increased to ensure
615 that the optimizer finds the global maximum of the Likelihood function. This can
616 be changed under "Optimization" in Toolbar.
617 o After calculation, the page 'MLE fit' is visited to present maximized results.
618
619 - **'Continuous LR (Integrated Likelihood based)'** (case **Weight-of-Evidence** and **'Database**
620 **search'**):
621
622 o Instead of optimizing the Likelihood of the unknown parameters, a **multivariate**
623 **integration** over the unknown parameters are applied both under hypothesis Hp and Hd.
624 o The accuracy of the integral depends on the specified '**relative error requirement'** (see
625 vignette for details).
626 ▪ Can be changed under "Integration" in Toolbar. Default is 0.005.
627 o In the output (see Figure 14), also the relative error of the LR is given in brackets.
628 o The integral requires that an **upper boundary** for the parameters mu (amount of DNA)
629 and sigma (coefficient of variation) are specified. As default these are 20000 and 1,
630 respectively. These values may be changed under "Integration" in Toolbar. See vignette
631 for details.
632 o Calculates LR-values directly and avoids visiting the tab 'MLE fit'.
633 ▪ Case **Weight-of-Evidence**: A message with LR pops up after calculation (see
634 Figure 14).
635 ▪ Case **'Database search'**: Database search results are shown directly after
636 calculation (goes to tab 'Database search').
637 o 'Continuous LR (Integrated Likelihood based)' is not possible for multiple replicates
638 and large number of loci since it doesn't evaluate on log-scale. Use the Maximum
639 Likelihood based method instead if the other method goes wrong.
640



641
642 Figure 14: The figure shows the calculated Weight-of-Evidence based the Integrated Likelihood based
643 continuous LR for the specified model in Figure 13.
644
645
646
647
648
649

650     -   **'Qualitative LR (semi-continuous)'** (case **Weight-of-Evidence**)
651

652           o   Performs a semi-continuous procedure where the distribution of the 'allele drop-out
653                probability given number of observed alleles' are utilized to infer a "conservative" LR.
654                   ▪   The model is purely qualitative which means it is only based on allele-
655                       information.
656           o   Goes directly to page Qual. LR.
657

658     -   **'Generate sample'** (case **'Generate sample'**):
659

660           o   A dataset (evidence sample and contributing references) will be randomly simulated
661                under the specified model under "Model specification".
662           o   Reference profiles may be imported and selected as assumed known in the hypothesis.
663           o   Detection threshold, stutter ratio, probability of drop-in and drop-in peak height
664                hyperparam may all be used in the simulation (**fst** are not used).
665           o   The unknown contributor profiles under the hypothesis will be randomly generated
666                using the selected population frequencies.
667           o   The simulated peak heights of the evidence in the dataset are entirely based on the
668                continuous model for assumed values of the model-parameters (**mu,sigma,xi,mx**).
669                Default these are given as **mu**=1000, **sigma**=0.15, **xi**=0.1, **mx**=(C:1)/sum(C:1), where C
670                is number of contributors.
671           o   Goes directly to page Generate data.

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691  # 3. <u>MLE fit</u>: ('Continuous LR (Maximum Likelihood based)')

692



693
694
695     Figure 15: The figure shows the <u>MLE-fit</u> GUI page after doing **continuous LR (Maximum**
696   **Likelihood based)** calculation (maximizing the continuous model with respect to the unknown
697       parameters for each of the specified hypothesis in figure 13) for **Weight-of-Evidence**.
698
699
700
701
702
703
704

<div align="center">ESTIMATES UNDER Hd (and Hp for case: **Weight-of-Evidence**)</div>

- **Parameter estimates**:

  - param: The unknown parameters in the model (see vignette for more details).
    - mx_i: Mixture-proportion for contributor 'i'.
    - mu: Expected amount of DNA.
    - sigma: Coefficient of variation.
    - xi: Stutter ratio (fraction of peak height that are stutter).

  - MLE: The optimized[1] parameters in the model which attains a maximum point of the likelihood function.

  - Std.Err.: The standard error of the parameter estimates in the model (see vignette for details).

- **Maximum Likelihood value**:
  - log10lik and Lik: The ten-logged and the original value of the Likelihood value attained from the optimization[1].

- **Further Action**:

  - **MCMC simulation** (see Figure 16):
    - Performs 'Markov Chain Monte Carlo (MCMC) random walk Metropolis' samples under the desired hypothesis.
      - Uses the mode and the covariance matrix attained from the optimization. See vignette for details.
    - The **first column** in the output shows the estimated posterior distributions for each of the unknown parameters in the model.
    - The **second column** in the output monitors the parameter samples in the simulation.
    - After sampling, the **acceptance rate** of the sampler is printed out to the terminal.
      - Acceptance ratio = number of accepted samples divided by number of proposed samples.
      - Ideally the acceptance rate should be around 0.2 to ensure that the parameter space has been fully explored.
        - Tweak '**variance of randomizer**' under MCMC in toolbar to change the acceptance rate.
    - User may **change number of required samples** in the simulation under 'MCMC' in toolbar.
    - The **purpose** of the MCMC simulation is to use it as an **exploratory tool** to see:
      - That the optimizer has found the global maximum.
      - The shape of the posterior distribution of the parameters.

---

[1] This may be only a local maximum point, not the global maximum (i.e. the Maximum Likelihood Estimate). Increase **number of start points** under "Optimization" in Toolbar to ensure a global maximum.

748
749 Figure 16: The figure shows the posterior density of the unknown parameters (first column) and
750 corresponding iteration values (second column) from the MCMC method under the hypothesis Hp:
751 "Suspect+1 unknown individual contributes to evidence evid1". The acceptance ratio was given as
752 0.35.
753

- **Deconvolution**:
  - Performs "Deconvolution" under the desired hypothesis. (See <u>Deconvolution</u> <u>(page 5)</u> for details.

- **Model validation** (Figure 17):
  - Uses a statistical hypothesis test to reject whether the maximum likelihood fitted model fits the observed peak heights (i.e. whether the gamma model assumption is reasonable).
  - Estimates the cumulative probability of the observed peak heights conditional on the other peak heights (see vignette for more details).

764
765
766
767
768
769

- Uses a one-sample Kolmogorov-Smirnov test to test if the observed cumulative probability deviates significant from the uniform distribution.
  - P-value from the test is printed out to terminal.
  - A textbox is shown when the P-value is lower than the significance level 0.05 (i.e. rejection of assumption).

Figure 17: Left subplot shows the "**Model validation**" under hd with p-value 0.37. Right subplot is "**Model validation**" under hp with p-value 0.29.

773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793

## WEIGHT-OF-EVIDENCE (case Weight-of-Evidence)

- Description:

  - The Weight-of-Evidence value is the ratio between the likelihoods of the two specified hypotheses Hp and Hd as specified in "Model specification".

  - The Weight-of-Evidence value is based on the continuous model as described in the vignette and handles allele drop-in, drop-out and stutter.

- **Join LR:**

  - LR: 'Likelihood value under optimization under Hp' divided by 'Likelihood value under optimization under Hd'
  - log10: The ten-logged value of LR.

794  -  **LR for each loci:**

795

796    o  The LR for each loci separately (given the parameter-modes under Hp and Hd). See
797       vignette for details.
798    o  Note: This will not be shown for evaluation of more than 30 loci

799

800

801                              FURTHER EVALUATION
802  -  **Optimize model more**:

803

804    o  The optimization procedure can be run again with the same specifications as selected in
805       "Model specification" to ensure that a global maximum is attained.
806    o  It is recommended to do this and check that the optimized Likelihood value is not
807       increased further.

808

809  -  **Database search (case: 'Database search'):**

810

811    o  A database search with the specified continuous model will be applied. (See <u>Database</u>
812       <u>search </u>for details.

813

814  -  **'Continuous LR (Integrated Likelihood based)' (case Weight-of-Evidence)**

815

816    o  See CALCULATIONS under section "<u>Model specification</u>".

817

818  -  **'Simulate LR distribution' (case Weight-of-Evidence)**

819

820    o  MCMC simulation will be applied both under Hp and Hd to provide a plot of a
821       "Bayesian" distribution of the LR where the uncertainty of the parameters in the
822       continuous model under both Hp and Hd are taken into account (see Figure 18).
823       ▪  Number of samples can be changed with **Set number of samples** under MCMC
824          in Toolbar (default is 10000 samples).

**Distribution of LR over posterior space of parameters**



825
826    Figure 18: The plot shows the distributed LR where the *a posteriori* density of the parameters in the
827    continuous model under both Hp and Hd are taken into account. *a posteriori* density are simulated
828    using the **MCMC simulation** (Figure 16 shows only Hp)**.**
829
830
831
832                                    SAVE RESULTS TO FILE
833
834    -   **'All results'**:
835           o   The parameter estimates with corresponding standard deviation errors estimates and the
836               likelihood values will be printed to file for all hypotheses on page (see below).

```
-------Estimates under Hd---------

param-MLE-Std.Err.
mx1-0.87124-0.06018
mx2-0.12876-0.06018
mu-1226.3- 116.6
sigma-0.42447-0.06491
xi-0.50195-0.07972

log10Lik=-111.3
Lik=5.024e-112

-------Estimates under Hp---------

param-MLE-Std.Err.
mx1-0.2296-0.0371
mx2-0.7704-0.0371
mu-1226.65-  93.23
sigma-0.33957-0.04491
xi-0.10902-0.06756

log10Lik=-107.3
Lik=5.36e-108
```

- **'Only LR results'**: **(case Weight-of-Evidence)**
  o The LR calculated values shown in WEIGHT-OF-EVIDENCE will be printed to file
    (see below).

| Marker | LR | log10LR |
|---|---|---|
| D3S1358 | 2.245e+00 | 0.35113 |
| VWA | 4.607e+00 | 0.66345 |
| D16S539 | 9.449e+00 | 0.97536 |
| D2S1338 | 6.980e+00 | 0.84384 |
| D8S1179 | 3.310e+01 | 1.51979 |
| D21S11 | 7.064e-01 | -0.15094 |
| D18S51 | 9.490e-02 | -1.02273 |
| D19S433 | 2.023e+00 | 0.30610 |
| TH01 | 3.843e+00 | 0.58467 |
| FGA | 9.067e-01 | -0.04253 |
| JointMLE | 1.067e+04 | 4.02814 |

# 4. Deconvolution:

853

854



855
856 Figure 19: The figure shows the Model Specification GUI page for doing **Deconvolution**. We
857 condition on the suspect, and assume one unknown in the hypothesis. Our model assumes
858 unknown "(n-1)- stutter" ratio, no allele drop-in and no theta-correction.
859
860
861 - Description:
862
863     o Deconvolution is applied for a specific hypothesis Hd as shown in Figure 19.
864     o The deconvolution conditions on the optimized parameters (i.e. the MLE fit in Figure
865         20) for the continuous model.
866     o The deconvolution result shows (see Figure 21) a ranked list of the **posterior**
867         **probabilities** of the combined genotype-profiles (see vignette for details).
868     o Since the deconvolution is based on the continuous model it may handle multiple
869         replicates, allele drop-in, drop-out and stutter.
870
871 - **Table**:
872
873     o The columns in the table (see Figure 21) show the resolved genotype for each
874         contributor in the specified hypothesis (per locus).

o The combined profiles are ranked due to their **posterior probabilities**.

876 o The ranked elements in the table ensures that the sum of the **posterior probabilities** are
877 at least 0.9999.
878 ▪ Can be changed under 'Deconvolution' in toolbar.
879 o Maximum length of table is default 10000.
880 ▪ Can be changed under 'Deconvolution' in toolbar.
881 o Note:
882 ▪ Having only sub-optimized parameters (in the **MLE fit**)will not give the most
883 likely genotypes.
884 ▪ Q-assignation is recommended to use since dropped out alleles are equally
885 threated and assigned as "99".

887 - **Save table:**

889 o The **full** table will be exported to a tabulator-separated text-file.



Figure 20: The figure shows the optimized parameters (i.e. the <u>MLE fit)</u> for the continuous model. The fitted model has the same "Further Action" possibilities as for "Weight-of-Evidence" and "Database search".

Figure 21: The figure shows the ranked table of deconvoluted genotype profiles for the unknown major contributor, when conditioning on the suspect profile. The table is ranked with respect to the posterior probability of different combined genotype profiles. The top ranked combined genotype profile is an outlier from the others which indicates that it is possible to extract the unknown profile (from figure 9 we see that this is a correct extraction).

# 5. Underline: Database search:

928

929



Figure 22: The figure shows the GUI page of the model specification for doing database search on the database file "databaseESX17". Our model assumes no "(n-1)-stutter", no allele drop-in and no theta-correction.

934

935

936  - Description:

937

938        o  The 'Database search' is very similar as the Weight-of-Evidence (see Figure 22) with
939           the only difference in that each individual in the reference-database is assumed as a
940           contributor in the hypothesis Hp. For each individual 'j' in reference-database we
941           calculate a LR-value LRj.

942       o   The user may choose between using peak heights in a 'Continuous LR' (**Maximum**
943         **Likelihood based** or **Integrated Likelihood based**)' calculation or ignoring the peak
944         heights in a 'Qualitative LR' calculation.
945
946     -   When selecting 'Continuous LR':
947
948       o   'Qualitative LR' is always calculated along with the 'Continuous LR' values.
949       o   The qualitative model assumes an allele drop-out parameter which is estimated.
950       o   The allele drop-in parameter in the qualitative model is set as default 0.05, but can be
951         changed with "**Set drop-in probability for qualitative model**" under 'Database search'
952         in the Toolbar.
953       o   No theta-correction is assumed in the qualitative model.
954       o   If "Continuous LR (Maximum Likelihood based)" calculation is used, the optimized
955         parameters under the Hd -hypothesis are first shown (see Figure 23).
956



957
958     Figure 23: The figure shows the optimized parameters (i.e. the <u>MLE fit)</u> for the continuous model
959   under Hd (with specifications as given in Figure 22). The fitted model has the same "Further Action"

possibilities as for "Weight-of-Evidence" and "Deconvolution". The user must push "**Database search**" for doing the actual database searching.

- When selecting 'Qualitative LR':

   o The "**Set drop-in probability for qualitative model**" under 'Database search' in the Toolbar is ignored.
   o The qualitative model assumes an allele drop-out parameter which is estimated.
   o The 'Continuous LR' calculation is ignored.

- Note:

   o The 'Continuous LR' calculation is based on the **continuous model** as given in the vignette and hence may handle allele drop-in, drop-out and stutter.
   o Continuous LR (Integrated Likelihood based) is not possible to use for replicates.
   o The reason for showing the MLE fitted parameters under Hd (see Figure 23) for "Continuous LR (Maximum Likelihood based)" calculation is that the user should have the possibility to check if the parameter estimates under Hd seems reasonable so he can go back and change the model specification.

- **Table** (see Figure 24):

   o '**Reference name'** is name of individuals given in the reference-database.

   o The table shows the ranked individuals in the database due to the continuous LR values (**contLR**), qualitative LR values (**qualLR**), number of matching alleles (**MAC**) or number of evaluating loci (**nLocs**).

   o **qual.LR** (Qualitative LR (semi-continuous model))
      ▪ Parameter for dropout probability is based on the median of 2000 samples from the 'distribution of dropout-probability'.
         • Number of required samples may be changed under 'Qual LR' in toolbar.
      ▪ For multiple evidences, the mean of the median is used as the dropout probability parameter.
      ▪ Assumes drop-in probability 0.05 as default. Can be changed under 'Database search' in toolbar.
      ▪ Assumes no theta-correction.

   o **MAC** (Matching allele counter) is number of alleles in the reference-profile which matches the evidence.
      ▪ Note: MAC is summed over the considered evidences.

   o **nLocs** is number of loci in the reference-profile which are used to calculate the contLR, qualLR and MAC.
      ▪ Note: Some references in the database may be missing loci which are presented in the evaluated evidence.

o Note:
1007             ▪ Maximum number of elements to view a 'Database search' result table is 10000.
1008                 This can be changed under 'Database search' in toolbar.
1009             ▪ Putting fst>0 may be very time-consuming since we require that individual 'j' is
1010                 a known non-contributor under Hd, and hence Hd is calculated for each
1011                 individual in database.
1012             ▪ If no allele drop-in is assumed under the continuous model, **cont.LR** is not
1013                 calculated for the non-fitting individuals in the database.
1014
1015     - **Save table:**
1016
1017         o The full table will be exported to a tabulator-separated text-file.
1018



1019
1020     Figure 24: The figure shows the table from the database search with specifications as given in Figure
1021     22 based on '**Continuous LR' (Maximum Likelihood based)"** calculations. The references are sorted
1022         due to the qualitative LR's (which assumes allele drop-out probability 0.08 and allele drop-in
1023                                     probability 0.05).
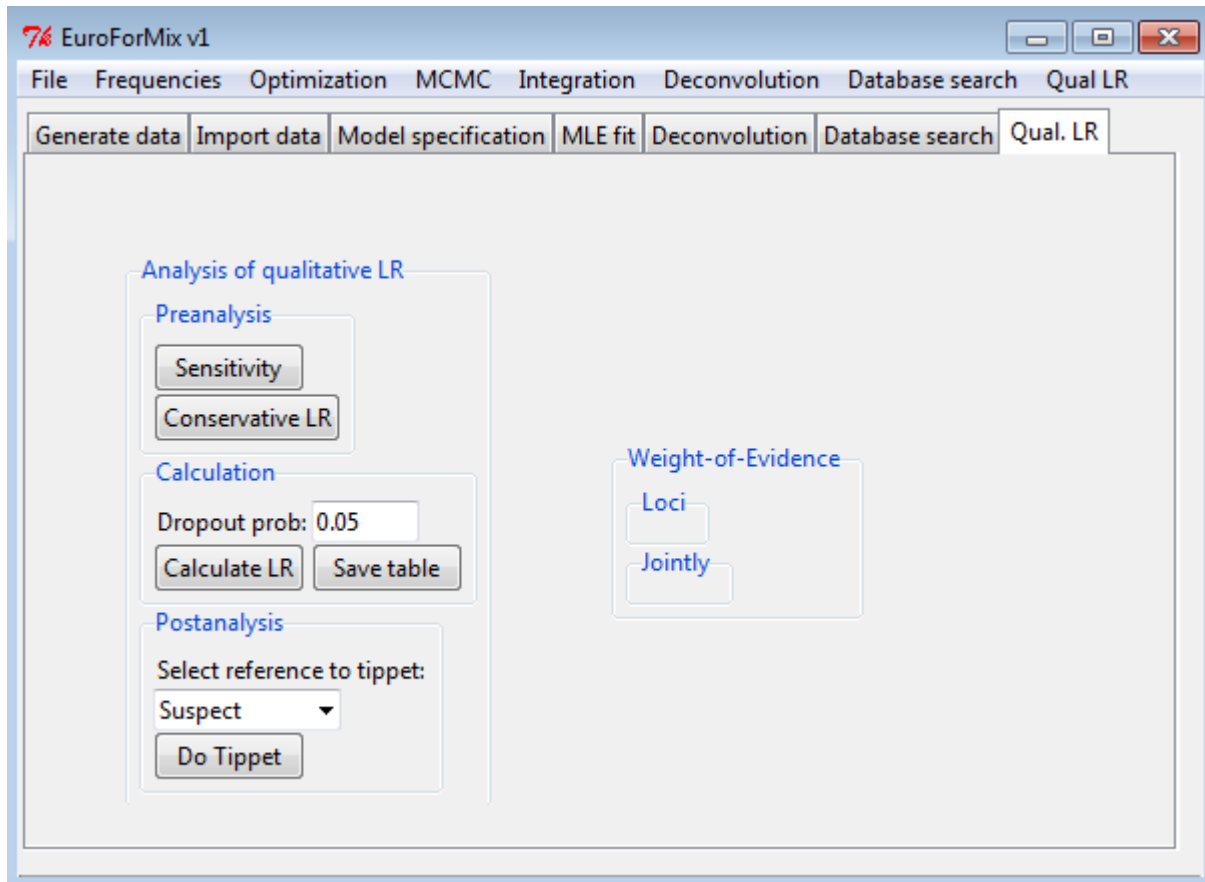
# 6. Qual. LR:



Figure 25: The figure shows the GUI page where the weight-of-evidence evaluation based on the qualitative model is done.

- Description:

  o This module samples from the distribution of the '*allele drop-out probability given number of observed alleles*' to evaluate the qualitative LR automatically. Also sensitivity plot as a function of allele-dropout probability and random man tippet analysis is implemented.


PREANALYSIS

- **Sensitivity:**

  o Plots the log10LR as a function of allele-dropout probability (see Figure 26).
    - The upper probability range and number of ticks can be changed under 'Qual LR' in the toolbar.
  o Note:
    - Lower range in sensitivity is 1e-6 (something small).

**Sensitivity plot**

1047
1048    Figure 26: The figure shows the plot of Weight-of-evidence (Likelihood Ratio) as a function of allele
1049                                drop-out probability.
1050
1051    -    **Conservative LR:**
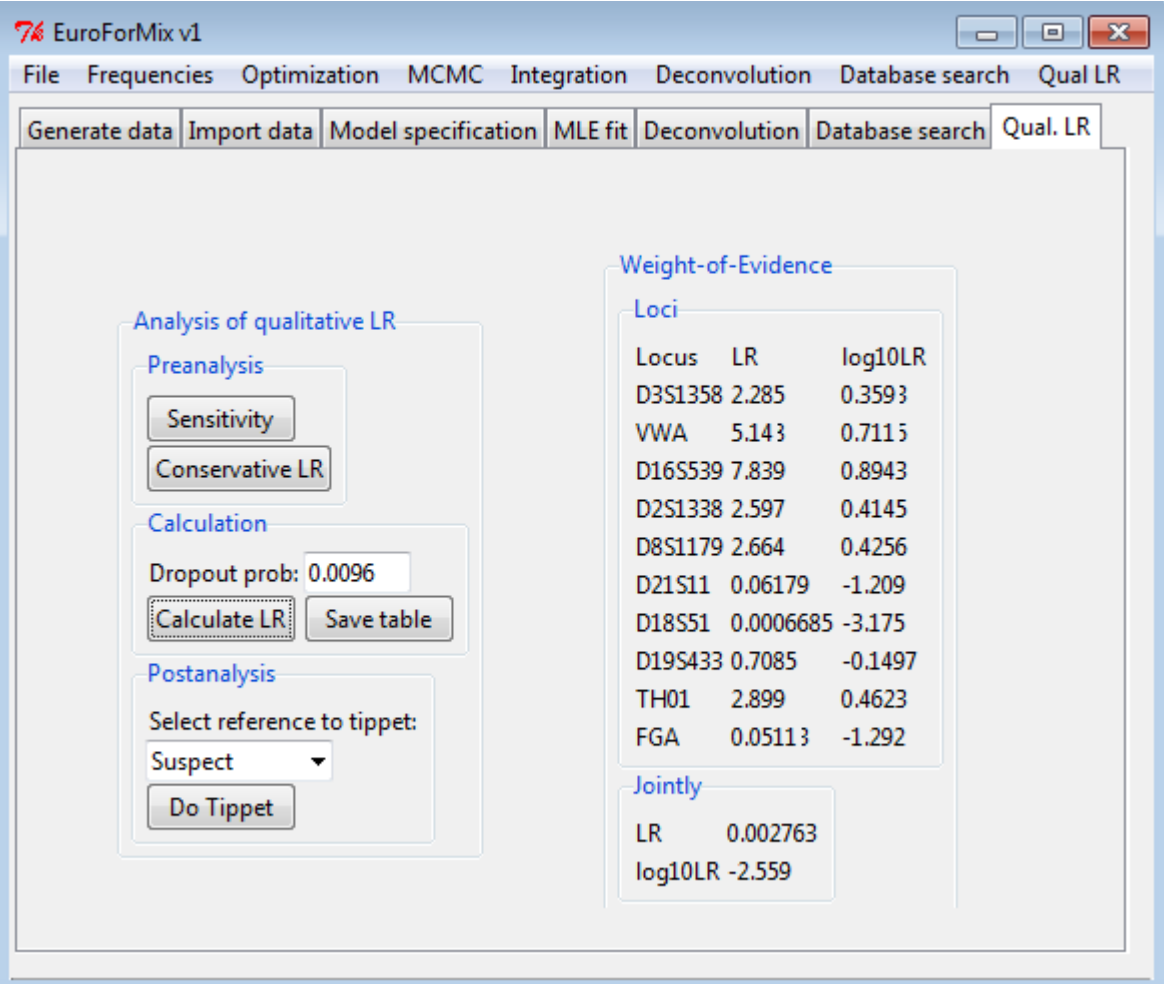1052
1053        o    By sampling from the "*allele drop-out probability given number of observed alleles in*
1054             *the evidence*"- distribution for the hypothesis Hp and Hd, the most 'conservative' LR
1055             (i.e. smallest) is automatically calculated and printed (see Figure 27 and Figure 28).
1056             ▪    The most "conservative" LR is found by following:
1057                  •    Take out the "alpha" and "1-alpha"-quantiles from the simulated 'allele-
1058                       dropout probability distribution' under both Hp and Hd.
1059                  •    The quantile (under both Hp and Hd) which gives the lowest LR is the
1060                       "conservative LR".
1061             ▪    The significance level "alpha" is given 0.05 as default.
1062                  •    This can be changed under 'Qual LR' in the toolbar.
1063             ▪    The number of required samples from the 'allele-dropout probability
1064                  distribution' is given 2000 as default.
1065                  •    This can be changed under 'Qual LR' in the toolbar.
1066             ▪    Note: If no samples are accepted from the allele-dropout probability
1067                  distribution', an error-message is provided to the user.
1068
1069        o    When more evidence samples are imported, the most 'conservative LR' over all samples
1070             is considered.
1071             ▪    The dropout probability quantiles are estimated for each of the evidence samples.
1072

1073
1074

```
[1] "For evidence evid1:"
[1] "Estimating quantiles from allele dropout distribution under Hp..."
         5%        95%
0.01089928 0.23859512
[1] "Estimating quantiles from allele dropout distribution under Hd..."
         5%         95%
0.009575226 0.223586744
         5%   95%
qqhp 0.0110 0.24
qqhd 0.0096 0.22
```

1075
1076    Figure 27: The plot shows the sampled 5% and 95% quantiles of the distribution of the '*allele drop-out*
1077                        *probability given number of observed alleles*'.
1078



1079
1080        Figure 28: The plot shows the conservative Weight-of-Evidence values (Likelihood Ratios) after
1081        pushing "**Conservative LR**". The most conservative estimated allele drop-out probability-quantile
1082    from Figure 27 was the 5% quantile under  Hd which gave 0.0096. Hence the table in this plot shows
1083                                the LR inserted for this value.
1084
1085
1086

CALCULATION

- **Dropout prob:**

  o The user may specify the assumed number of allele dropout-probability.

- **Calculate LR**

  o Instantly calculates the LR for the given user-specified allele dropout probability in "**Dropout prob**".

- **Save table:**

  o Saves the weight-of-evidence calculated LR results to a selected file.

POSTANALYSIS

- **Selection of reference to tippet:**

  o A drop-down list of references which are conditioned under Hp but not under Hd.

- **Do Tippet:**

  o Random tippet samples are provided by replacing the selected reference (under the drop-down list in the hypothesis Hp) with a random individual from the population and then calculate his LR. A vast amount (default is 1e6) of random tippets is simulated to determine the tippet-distribution.
    ▪ The mean, standard errors of LR and log10LR-quantiles (1%, 5%, 50%, 95%, 99%) are printed out to terminal (see Figure 29).
    ▪ A plot of the cumulative distribution of log10LR will be shown (see Figure 30).
    ▪ Number of tippets can be changed under 'Qual LR' in the toolbar.
  o If weight-of-evidence has been calculated:
    ▪ The reporting LR for the "tipped individual" is superimposed as a blue line to the plot (see Figure 30).
    ▪ The discriminatory metric (log10LR-q99%) is printed out to terminal (see Figure 29).
  o Note: Precalculations are always done previous to the tippet-sampling, therefore the number of tippets are only limited to make the plot.

```
[1] "Precalculating for tippet plot..."
[1] "Simulating 1e+06 tippets..."
[1] "Mean of samples = 0.703689922273713"
[1] "Standard Error of samples = 0.677592812037152"
           1%          5%          50%          95%          99%
    -32.630258  -28.495191  -18.700169   -9.622545   -6.242242
[1] "Discriminatory metric (log10(LR) - q99) = 3.6836388776166"
```

Figure 29: The plot shows the printed tippet information to the terminal when replacing the "Suspect" in hypothesis Hp with a random man (a tippet). Number of tippets simulated, mean and standard errors of LR and log10LR-quantiles (1%, 5%, 50%, 95%, 99%) are printed out to terminal (see Figure 29). Also the discriminatory metric, the distance between the observed log10LR for the suspect and log10LR-99%-tippet-quantile is given.
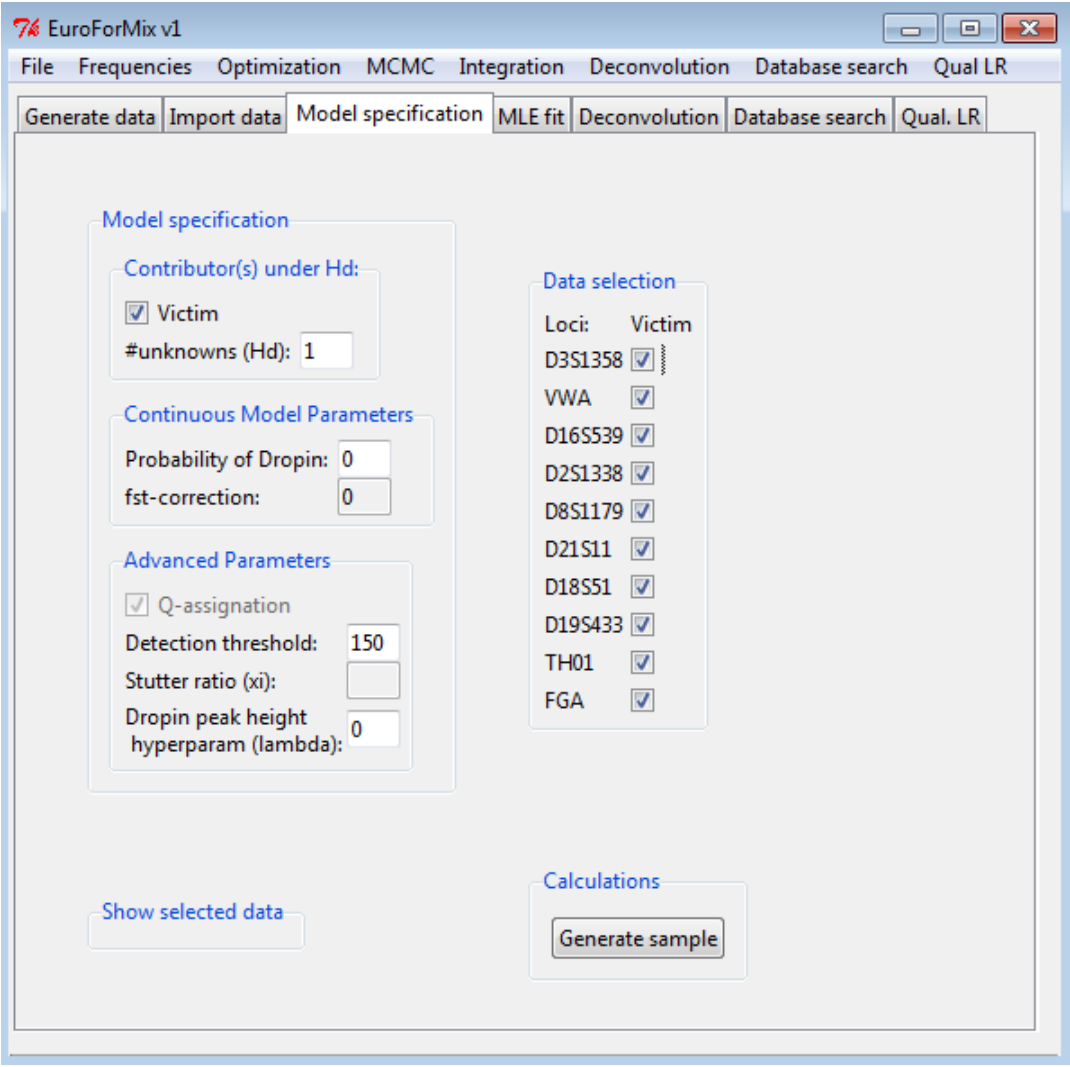


Figure 30: The figure shows a cumulative distribution of 1000000 log10LR tippets, where each tippet is based on replacing the "Suspect" in hypothesis Hp with a random man from the population. The reporting LR for the "tipped individual" (i.e. "Suspect in this case) is superimposed as a blue line to the plot.

# 7. Generate data:

Figure 31: The figure shows the Model specification GUI page for generating allele with corresponding peak heights from the continuous model for a given specified model. From here we will generate data which are contributed from a known Victim profile and an unknown individual. We assume a detection threshold of 150 rfu and no allele drop-in is considered.

- Description:

   o Generates alleles using the population frequencies and simulates peak heights for a specified hypothesis (see Figure 31) using the continuous model.
   o The generation may simulate allele-dropout, drop-in (with a peak height model) and stutter (see Figure 32).
      ▪ Allele-dropout is indirectly simulated by falling below the defined threshold.

1156
1157 Figure 32: The figure shows the <u>Generate data</u> GUI page which shows the generated alleles and
1158 corresponding peak heights (under **Evidence**) for the given selected set of parameters under
1159 **Parameters**. The true contributors are given under **Reference(s)**.
1160
1161
1162
1163
1164

1165

1166    -    **Parameters**:

1167

1168              o   **mu**:  amount of DNA
1169              o   **sigma:** coefficient of variance
1170              o   **xi:** stutter ratio
1171              o   **mx=(mx1,…, mxC):** mixture proportion for contributor 1,..,C.
1172                        ▪   Note: **mx** will be normalized if it's not already.

1173

1174    -    **Edit**:

1175

1176              o   **Loci**: Loci name of the population frequency used to generate the dataset.
1177              o   **Evidence**: The allele information is given in the left column while the peak height
1178                  information is given in the right column. Each element **needs to be** separated with ",".
1179              o   **Reference**: The alleles of the true contributors to the generate evidence is sequentially
1180                  shown in each column.
1181              o   All the loci names, evidence-allele and heights and reference-alleles may be edited
1182                  before storing (See Figure 32).

1183

1184    -    **Import/Export**:

1185

1186              o   **Save data**:
1187                        ▪   Stores the generated (and possible edited) evidence- or reference-profile to a file.
1188                        ▪   Extension .csv added automatically.

1189

1190              o   **Load data:**
1191                        ▪   Loads profiles from file into the selected entries (evidence or reference).
1192                                •   This is useful for generating random evidence samples where loaded
1193                                    references are conditioned on.
1194                        ▪   Note:
1195                                •   If any locus is missing from the loaded evidence or reference file, the
1196                                    edit-cell will be empty.
1197                                •   The order of the loci in the file does not matter.

1198

1199    -    **Further action:**
1200              o   Generate again**:** Make a new simulation of the evidence sample using the selected
1201                  values of the parameters under **Parameters**.
1202              o   Plot EPG: Plots the generated (and possible edited) evidence in a EPG-plot.
1203                        ▪   It will use the "kit" selected under "Import Data"-page.
1204                        ▪   See ?plotEPG to see which kit-formats that are supported in the EPG.

1205

1206

1207

1208

1209

1210

# (C) To be implemented in a future version:

- Label the alleles of the selected references to the EPG-plot.
- Warning if exp(lik)=0 when lik>-Inf (happens for INT calculations)
- Empty loci will not be removed when imported to the software. They will be considered as a full dropped out loci in the evaluation.