1    **High-quality, genome-wide SNP genotypic data for pedigreed germplasm of the diploid outbreeding**

2    **species apple, peach, and sweet cherry through a common workflow**

3

4    Stijn Vanderzande[1,‡,*], Nicholas P Howard[2,3,‡], Lichun Cai[4], Cassia Da Silva Linge[5], Laima Antanaviciute[5],

5    Marco CAM Bink[6,7], Johannes W Kruisselbrink[6], Nahla Bassil[8], Ksenija Gasic[5], Amy Iezzoni[4], Eric Van de

6    Weg[9], Cameron Peace[1]

7

8    1    Department of Horticulture, Washington State University, PO Box 646414, Pullman, WA 99164,

9        USA

10   2    Department of Horticultural Science, University of Minnesota, St Paul, MN 55104, USA

11   3    Institute of Biology and Environmental Sciences, Carl von Ossietzky Universität, 26129

12       Oldenburg, Germany

13   4    Department of Horticulture, Michigan State University, East Lansing, MI 48824, USA

14   5    Department of Plant and Environmental Sciences, Clemson University, Clemson, South

15       Carolina, USA

16   6    Biometris, Wageningen UR, PO Box 16, 6700AA, Wageningen, The Netherlands.

17   7    Research & Technology Center, Hendrix Genetics, P.O. Box 114, 5830 AC Boxmeer, The

18       Netherlands

19   8    USDA-ARS, National Clonal Germplasm Repository, 33447 Peoria Road Corvallis, OR 97333,

20       United States of America

21   9    Plant Breeding, Wageningen UR, PO Box 386, 6700AJ, Wageningen, The Netherlands

22   ‡    These authors contributed equally

23   *    Corresponding author: stijn.vanderzande@wsu.edu

24

**Abstract**

High-quality genotypic data is a requirement for many genetic analyses. For any crop, errors in genotype calls, phasing of markers, linkage maps, pedigree records, and unnoticed variation in ploidy levels can lead to spurious marker-locus-trait associations and incorrect origin assignment of alleles to individuals. High-throughput genotyping requires automated scoring, as manual inspection of thousands of scored loci is too time-consuming. However, automated SNP scoring can result in errors that should be corrected to ensure recorded genotypic data are accurate and thereby ensure confidence in downstream genetic analyses. To enable quick identification of errors in a large genotypic data set, we have developed a comprehensive workflow. This multiple-step workflow is based on inheritance principles and on removal of markers and individuals that do not follow these principles, as demonstrated here for apple, peach, and sweet cherry. Genotypic data was obtained on pedigreed germplasm using 6-9K SNP arrays for each crop and a subset of well-performing SNPs was created using ASSIsT. Use of correct (and corrected) pedigree records readily identified violations of simple inheritance principles in the genotypic data, streamlined with FlexQTL™ software. Retained SNPs were grouped into haploblocks to increase the information content of single alleles and reduce computational power needed in downstream genetic analyses. Haploblock borders were defined by recombination locations detected in ancestral generations of cultivars and selections. Another round of inheritance-checking was conducted, for haploblock alleles (i.e., haplotypes). High-quality genotypic data sets were created using this workflow for pedigreed collections representing the U.S. breeding germplasm of apple, peach, and sweet cherry evaluated within the RosBREED project. These data sets contain 3855, 4005, and 1617 SNPs spread over 932, 103, and 196 haploblocks in apple, peach, and sweet cherry, respectively. The highly curated phased SNP and haplotype data sets, as well as the raw iScan data, of germplasm in the apple, peach, and sweet cherry Crop Reference Sets is available through the Genome Database for Rosaceae.

2

49

50  **Introduction**

51  A high-quality, mostly error-free genotypic data set is imperative to obtain reliable results in many

52  downstream genetic analyses. The results of genetic analyses can be influenced by even low rates of

53  genotyping errors [1]. For example, the size of genetic maps and order of markers therein are affected

54  by errors in genotypic data [2–4]. Inaccurate genotypic data will also lower the power, accuracy, and

55  resolution of linkage studies and increase the number of false marker-locus-trait associations [5–7]. The

56  number of observed (double) recombinants is inflated by errors in genotypic data [8]. Incorrect calling of

57  recombinations in turn leads to incorrect determination of haploblock limits and assignment of

58  haplotypes [9]. Finally, incorrect genotype calls can lead to incorrect imputations of missing data or even

59  the improper adjustment of correct data to ensure the data is consistent with Mendelian inheritance

60  [10].

61

62  There are several reasons for the occurrence of errors in a genotypic data set. Incorrect information

63  about a sample's identity, e.g., due to mixing up or mislabeling samples, causes an individual to be

64  matched with the wrong data [1]. In clonally propagated crops, mislabeling errors can easily spread

65  when individuals that are not true-to-type are used as parents or as base plants to create new

66  propagules. Available pedigree information for an individual can be incorrect, causing incorrect

67  enforcement of allele assignments. In fruit cultivars, numerous pedigree records have been confirmed or

68  updated with the help of genetic markers [11–23]. Biological reasons such as unexpected mutation,

69  insertions or deletions in the DNA sequence containing markers, and gene conversion can lead to

70  inconsistencies in genotype calls and propagate errors through the data set [1]. Technician errors can

71  also introduce errors in a data set, such as when lab protocols are not applied correctly (Hoffman and

3

72    Amos 2005) or when multiple large data sets with disparate formats are integrated and edited. Finally,

73    technological and software limitations and failures can also lead to the presence of errors [1].

74

75    SNPs have become the genetic marker of choice for many genetic analyses but, with their increased use

76    and increasingly large numbers that can be generated, manual data curation has become more

77    challenging. SNPs are ubiquitous within the genome and allow for simultaneous screening of many

78    thousands of polymorphic loci via SNP arrays, Genotyping-By-Sequencing, or resequencing [24,25]. SNP

79    arrays provide consistent information between individuals and have been developed for clonally

80    propagated crops, such as the 8K apple array [26], 9K peach array [27], and 6K cherry array [28]

81    developed by international teams led by RosBREED; the GrapeReSeq 18K Vitis array [29]; the 20K apple

82    array developed by FruitBreedomics [30], all on the Illumina Infinium® platform, and the strawberry 90K

83    Axiom array [31], and the 480K apple array by FruitBreedomics on the Affymetrix axiom platform [32].

84    Genotyping each individual relies on the automated scoring of thousands of SNPs. As thousands to

85    millions of SNPs are being assessed on a large set of individuals, even a low error rate in SNP scoring can

86    correspond to a high absolute number of errors. As the number of SNPs on an array increase, it becomes

87    more time-consuming and less feasible to manually review all automated SNP calls to identify potential

88    errors.

89

90    For SNP arrays, incorrect genotype assignment using automated SNP scoring software occurs when

91    intensity plots deviate from expected patterns. Automated genotyping is based on the association of

92    specific alleles to different fluorescent molecules, the detection of these fluorescent molecules, the

93    clustering of individual-marker data points according to intensity ratios between the different

94    fluorescent dyes across multiple individuals into distinct regions of a genotype-calling space, and the

95    final assignment of these clusters to genotypes. Examples of deviations that are observed in the

96    intensity plots are the presence of additional clusters or clusters that have shifted from their expected

97    location in the intensity plot. The presence of additional clusters or shifted clusters can be attributed to

98    additional regions that bind to the SNP's probe [33]. Sequence similarity of these regions with the

99    intended target is caused by either local sequence repetition or presence of paralogous regions in the

100   genome. The presence of these highly similar sequences can lead to multi-locus segregating SNP

101   markers that cannot be adequately called. The calling of a single segregating locus might also be

102   hampered by the background signal of targeted but non-segregating gene copies (ASSIsT Reference

103   Manual p17 [34]). The presence of one or more additional SNPs, insertions, or deletions in the probe-

104   binding region can lead to reduced or loss of binding affinity for the SNP's probe and thereby to the

105   presence of additional clusters, both of which can lead to incorrect genotype scoring of some SNPs [33].

106

107   No systematic workflow exists to efficiently detect and resolve all types of errors from a genotypic data

108   set for pedigreed germplasm. Methods and software exist to tackle specific types of errors. For example,

109   the ASSIsT software was developed for use with Illumina Infinium® arrays to identify which SNPs show

110   robust results, which SNPs might have genotype calling errors due to alleles with reduced affinity or null

111   alleles, and which SNPs are monomorphic or failed completely [35]. Another example is the aggregation

112   of linked SNPs into a single genetic locus, called haploblock, which facilitates tracking the inheritance of

113   alleles within a pedigree and subsequent identification of inheritance inconsistencies [36]. Despite the

114   existence these and other methods and software, an effective way to combine these methods has not

115   been described.

116

117   Here we describe a curation workflow for high-resolution genetic marker data that identifies and

118   resolves errors to obtain a robust set of genotypic data. The workflow maximizes the genotypic data

119   obtained from high-throughput genome-scanning tools while minimizing the time needed to identify

5

120    and remove errors. The workflow resulted from curation needs in the multi state and multi-crop USDA-

121    SCRI project RosBREED [37–39] and the European project FruitBreedomics [40–42]. The workflow is

122    demonstrated for three tree fruit crops, apple, peach, and sweet cherry, using the RosBREED germplasm

123    sets [43]. The resulting genotypic data sets can be used by researchers to reconstruct pedigrees,

124    establish quantitative genetic relationships, identify and validate quantitative trait loci (QTLs), and trace

125    allele sources, leading to valuable practical and scientific genetic insights – with high confidence in the

126    obtained results.

127

128    **Material and Methods**

129

130    *Plant material*

131    The apple, peach, and sweet cherry collections used in this study, referred to as the 'Crop Reference

132    Sets', were created to represent U.S. breeding germplasm [43] for the RosBREED project [37]

133    (www.rosbreed.org) and consisted of 451, 426, and 269 individuals for apple, peach, and sweet cherry,

134    respectively (Tables S1-S3). Three apple breeding programs (Washington State University, the University

135    of Minnesota, and Cornell University), three peach breeding programs (University of Arkansas, Clemson

136    University, and Texas A&M University), and one sweet cherry program (Washington State University)

137    each contributed additional germplasm to complement the Crop Reference Sets and better represent

138    their important breeding parents [43]. These additional 'Breeding Pedigree Sets' consisted of 172, 139,

139    and 167 apple individuals, 117, 289, and 143, peach individuals, and 259 sweet cherry individuals,

140    respectively. The sweet cherry Breeding Pedigree Set was later made publicly available and became part

141    of the sweet cherry Crop Reference Set. Genotypic data of the other Breeding Pedigree Sets were

142    included as part of the data curation but individual identities of this private germplasm are not provided.

143

144    To reduce the trimming of pedigrees (as described under 'Haploblock and haplotype generation' below),

145    the genotype calls of 18 additional apple individuals genotyped with the 20K SNP array in the

146    FruitBreedomics project [42] or genotyped with the 8K SNP array at KU Leuven, Belgium (Table S1) were

147    added to the data set to complete genotypic data of key ancestors.

148

149    *Initial parentage information*

150    Initial parentage information was collected as part of the germplasm creation as described by Peace and

151    co-workers (2014) [43]. For each breeding program, breeders provided pedigree records for their

152    seedlings, selections, and released cultivars. Other pedigree records were based on historical records

153    and available literature and were included for all progenitors, regardless of availability so that all

154    progenitors terminated in founders (individuals with two unknown parents).

155

156    *DNA extraction and iScan*

157    DNA extraction was conducted for apple, peach, and sweet cherry as described by Chagné and co-

158    workers (2012) [26], Verde and co-workers (2012) [27], and Peace and co-workers (2012) [28],

159    respectively. Genomic DNA from each individual was purified using the E-Z 96 Tissue DNA Kit (Omega

160    Bio-Tek, Inc., Norcross, GA, USA). DNA was quantitated with the Quant-iT™ PicoGreen® Assay

161    (Invitrogen, Carlsbad, CA, USA), using the Victor multiplate reader (Perkin Elmer Inc., San Jose, CA, USA).

162    DNA concentrations were adjusted to a minimum of 50 ng/µl, in 5 µl aliquots. For apple, DNA samples

163    were run on the Illumina Infinium® 8K apple SNP array [26] with iScans either at the Biotechnology

164    Platform of the Agricultural Research Council (Pretoria, South Africa) or at the Research Technology

165    Support Facility at Michigan State University (East Lansing, MI, USA), following the manufacturer's

166    protocol (Illumina Inc.). For peach and sweet cherry, DNA samples were run on the 9K peach SNP array

167    [27] and 6K cherry SNP array [28], respectively, with an iScan at the Research Technology Support

168     Facility at Michigan State University (East Lansing, MI, USA), following the manufacturer's protocol

169     (Illumina Inc.).

170

171     *Initial genetic maps*

172     For each crop, available genetic maps were used as a framework to determine the initial order of

173     reliable SNPs. Reliable SNPs (obtained as described under 'Subset of reliable SNP obtainment' below)

174     that were not present in available genetic maps were incorporated by comparing their physical positions

175     to those of flanking SNPs that were present in available genetic maps.

176

177     For apple, an integrated genetic map based on five full-sib families with 'Honeycrisp' as common parent

178     [20] was used as a framework to help align additional SNPs on the 8K array. The relative order of SNPs in

179     the map of Howard and co-workers (2017) [20] was adjusted to be consistent with the 'Golden

180     Delicious' double haploid genome sequence v1.1 [44] whenever this did not result in false detection of

181     double recombination for the original mapping populations. Then, SNPs that were included in the iGL

182     map [45] but not included by Howard and co-workers (2017) [20] were aligned based on relative marker

183     order between common markers of both maps and the 'Golden Delicious' double haploid genome

184     sequence v1.1 [44]. In cases of conflict between the iGL map and the reference genome, only the iGL

185     map was used as reference. Genetic positions of newly added SNPs were determined so that, in the new

186     map, they had the same position relative to the position of flanking markers as these SNPs did in the iGL

187     map. Finally, any remaining unmapped SNPs were positioned based solely on relative physical positions

188     according to the 'Golden Delicious' double haploid genome sequence v1.1 [44]. When the genetic

189     position in the iGL map was known for repositioned or newly added SNPs, their genetic position in the

190     new map was determined so that they had the same position relative to the position of flanking markers

191     as they did in the iGL map. When no genetic position in the iGL map was available, the genetic position

192     was determined so that, in the new map, they had the same position relative to the position of flanking

193     markers as they did in the physical genome. In peach, genetic positions were based on the peach

194     physical position of peach genome v2.0 [46]. The peach physical map was scaled to an approximate

195     genetic map by using a conversion factor where every 1 Mb corresponded to 4 cM. For sweet cherry,

196     genetic positions were determined by aligning and integrating the physical positions using peach

197     genome v2.0 [46] with the sweet cherry 'Regina' × 'Lapins' SNP linkage map [21,47].

198

199     *Workflow procedures*

200     Throughout the workflow, several software packages were used. Below are described the main

201     procedures used in the workflow, the associated software and parameter settings, and output files

202     used. The order in which each functionality was used in the workflow is reported in Results section

203     'Steps of the data curation workflow'.

204

205         *Initial genotypic data obtainment (GenomeStudio®)*

206     iScan output was converted to 'AA', 'AB', and 'BB' genotype calls for each SNP marker with the

207     Genotyping module of GenomeStudio® v2011.1 (Illumina Inc., San Diego, CA, USA) using a sample sheet

208     to load sample intensities and a 'Gen Call' Threshold of 0.15 to assign samples to a genotype cluster. The

209     sample sheet was adjusted in Microsoft Excel as follows before using it as input for GenomeStudio®:

210     • The sample sheet was saved as an 'xls(x)' file to avoid the loss of 'SentrixBarcode' information

211         that occasionally occurs when saving it as a '.csv' file.

212     • When individuals were separated over multiple iScan runs and sample sheets, the '[Data]'

213         sections of each sample sheet were combined into one.

214     • A copy of the 'Sample_ID' column in the '[Data]' section was added and named

215         'Sample_Original'.

9

216     •    Sample names in the 'Sample_ID' were adjusted to remove any spaces or special characters

217        (needed for some software) and avoid long names or names that could be interpreted as dates

218        (or other special formats) by Excel.

219     •    Duplicate and parental information was added to the 'Replicate', 'Parent1', and 'Parent2'

220        columns considering the adjusted names in the 'Sample_ID' column.

221     •    The resulting sample sheet was saved both as a '.xlsx' file for future editing and as a '.csv' file to

222        serve as an input file for GenomeStudio®.

223

224     *Low-quality and non-diploid sample identification (GenomeStudio® and R)*

225 Quality and ploidy were assessed using each sample's B-allele frequencies calculated by

226 GenomeStudio®. In GenomeStudio®, the histogram of the B-allele frequency was plotted for each

227 individual by opening the 'Histogram plot' function of the 'Full Data Table', choosing the first individual

228 in the 'Columns' section, and then choosing 'B Allele Freq' in the 'Sub Columns' section. The histogram

229 for the 'B-allele frequency' could then be plotted for each individual by scrolling through the individuals

230 in the 'Columns' section. Samples were considered of good quality when a clear heterozygous peak was

231 observed around 0.5 with almost no SNPs having a B-allele frequency between 0.125 and 0.375 and

232 between 0.625 and 0.75. In contrast, samples of poor quality showed no clear heterozygous peak

233 around 0.5 and had many SNPs with a B allele-frequency between 0.125 and 0.375, and between 0.625

234 and 0.75. Individuals that showed more than three peaks in the histogram were classified as polyploid.

235 Individuals that showed a 'shoulder' on the AB peak were classified as putative aneuploids and were

236 examined further in B-allele frequency plots according to Chagné and co-workers (2015) [48], below.

237

238 To create B-allele frequency plots according to Chagné and co-workers (2015) [48], a subset of SNPs was

239 created by applying the filter parameters described in Table S4A in the 'SNP Table' of GenomeStudio®.

240    Next, the 'Full Data Table' of GenomeStudio® was adjusted to only contain the B-allele frequency of

241    each sample: in the 'Column Chooser' function of GenomeStudio®, 'B Allele Freq' was added to the

242    'Displayed Subcolumns' section while all other subcolumns were removed from this section. The

243    resulting 'Full Data Table' was exported using the 'export displayed data to a file' function. The exported

244    'Data Table' was further adjusted to the following format: the first column contained the SNPs name,

245    the second column contained the SNP's cumulative position, and all subsequent columns contained the

246    samples' B-allele-frequencies.

247

248    Each SNP's cumulative genomic position was determined as follows: the chromosome number

249    corresponding to the SNP was multiplied by the power of ten which ensured that the outcome was

250    larger than any possible position within any chromosome (e.g., if the largest physical position within any

251    chromosome was 456,437 bp, all chromosome numbers were multiplied by 1,000,000 or $10^6$ as this is

252    the first power of 10 that is larger than 456,437. Similarly, if the largest genetic position within any

253    chromosome was 145 cM, each chromosome number was multiplied by 1000 or $10^3$). Then, the physical

254    or genetic position within the chromosome was added to the adjusted chromosome number to obtain

255    the cumulative genomic position of that SNP. The resulting file was then loaded into R [49].

256

257    An ad hoc R-script (Document S1) generated a pdf file that contained a plot for each individual where 'B-

258    allele frequency' values were plotted for the subset of SNP markers that were ordered according to their

259    cumulative position on a genetic linkage map or reference genome sequence. 'B-allele frequency' values

260    were expected to be 0, 0.5, or 1 for diploids. Diploid samples were considered of sufficient quality when

261    *almost no* SNPs (<0.3% of the subset) were observed between 0.125-0.375 and 0.625-0.875. In contrast,

262    a sample was considered of intermediate or poor quality when many SNP markers (0.3%-3% and >3%,

263    respectively) showed an intermediate or large discrepancy. For triploids, 'B-allele frequency' values were

11

264    expected to be 0, 0.33, 0.66, and 1 for all chromosomes while values of 0, 0.25, 0.5, 0.75, and 1.0 were

265    expected for tetraploids. Aneuploids had a diploid pattern for most chromosomes and a haploid or

266    polyploid pattern for others. Individuals classified as poor quality, polyploid, and aneuploid were

267    excluded from further analyses.

268

269    Samples were excluded from various input files and from the genotype clustering in GenomeStudio® by

270    choosing them in the 'Samples Table' and then choosing the 'Exclude Selected Samples'. SNPs were then

271    re-clustered by choosing the 'Cluster All SNPs' of the 'Analysis' section. All statistics were updated when

272    prompted.

273

274        *Subset of reliable SNP obtainment (ASSIsT)*

275    The 'Final report' and 'DNA report' input files were created as described in the ASSIsT Reference Manual

276    [34]. Briefly, a 'Final Report' and 'DNA Report' were generated using the 'Report Wizard' under the

277    'Reports' option of the 'Analysis' section. The best 'redo' was chosen based on the '10th Percentile GC

278    score' and excluded samples were removed from the report. For the 'Final Report', 'GTScore', 'Theta',

279    and 'R' were added to the default 'Displayed Fields' and data was grouped 'by SNP'. For the 'DNA

280    Report', samples were exported by 'Sample ID'.

281

282    The pedigree input file was created in Excel by copying the 'Sample_ID', 'Parent1', and 'Parent2'

283    columns from the '[Data]' section of the sample sheet used to create the GenomeStudio® project,

284    adjusting the column names to '//SampleID', 'Mother', and 'Father', respectively, and saving the

285    resulting file as a tab-delimited text file. The (optional) map was created in Excel by having the SNP

286    Names as given by GenomeStudio® in the first column and their corresponding chromosome and

287    position within the chromosome (either physical or genetic) as the second and third column,

12

288     respectively. Column names were set to '//SNPid', 'Chromosome', and 'Position' and the resulting file

289     was saved as a tab-delimited text file

290

291     All input files were loaded into ASSIsT v1.01 [35] using the 'Select' button. Then, parameters were set

292     using the 'Set' button as described in Table S4B depending on the 'Population type' used. ASSIsT

293     distinguished eight marker classes, which were re-grouped into the following five categories:

294     •    Robust SNPs: having less than 5% No Call Rate and all three possible clusters (AA, AB, and BB)

295            present in the germplasm set. In ASSIsT, these SNPs were classified as 'Robust',

296            'OneHomozygRare_HWE', 'OneHomozyRare_NotHWE', and 'DistortedAndUnexSegreg'

297     •    Two cluster SNPs: having less than 5% No Call Rate and one of the homozygous clusters (AA or

298            BB) absent in the germplasm set. In ASSIsT, these SNPs were classified as 'ShiftHomo'

299     •    Null-allele SNPs: having a probable null allele, classified as 'NullAllele-Failed' in ASSIsT

300     •    Monomorphic SNPs: having no polymorphism, as in ASSIsT

301     •    Failed SNPs: having more than 50% No Call Rate, poor clustering, or low intensity, as in ASSIsT

302

303     Results of SNP performance in ASSIsT were exported to the 'Summary' and 'Custom SNP information

304     table'. Genotype calls were saved in 'Custom gtypes' to be used in the R-script that checked pedigree

305     records (described below in 'Pedigree records verification'). PLINK input files were generated to check

306     for unknown duplicates within the data (described below in 'Duplicate individuals detection') and

307     FQ_DataPrepper input files were created to easily generate FlexQTL input files using FQDataPrepper

308     (described below in 'Genotyping error detection and adjustment'). Genotype calls for the 'Robust SNPs'

309     category were automatically reported in ASSIsT output files whereas other categories were considered

310     to contain failed SNPs and thus their genotype calls were not automatically reported. To include

13

311     genotype calls of the 'Two cluster SNPs', genotype calls of such SNPs were extracted from

312     GenomeStudio® and added to the data files manually.

313

314         *Duplicate individuals detection (GenomeStudio® and Plink)*

315     Genotypic data of known mutants and duplicates were compared to ensure their genotypic data were

316     matching using the 'Reproducibility and Heritability' report of GenomeStudio®

317     (Analysis>Reports>Reproducibility and Heritability Report>with Calculating Errors). The data set was also

318     screened for individuals with (unknown) identical genotypic data using Plink 1.9 [50] (https://www.cog-

319     genomics.org/plink2). Plink input files generated with ASSIsT were copied into the folder that contained

320     the PLINK executable (plink.exe). Then, a 'command window' or 'PowerShell window' was opened in this

321     folder and the 'plink.exe --file [*filename]* –missing-genotype - --genome full' or '\plink.exe --file

322     [*filename]* –missing-genotype - --genome full' command was given, respectively, where [filename] was

323     the name of the PLINK input files used. The resulting 'plink.genome' was opened in Excel and the

324     'PI_HAT' column was used to represent the proportion of identity-by-descent (IBD) between each pair of

325     individuals. Pairs of individuals with an IBD proportion higher than 97% were considered to be

326     duplicates because at this stage all known duplicates shared an IBD proportion of at least 97%. If

327     individuals were true duplicates, only one was kept in the data set. If pedigree records differed between

328     duplicate individuals, pedigree records were used to identify trueness-to-type as described below. True-

329     to type individuals were kept in the data set and individuals that were not true-to-type were targeted

330     for DNA re-sampling. Where two unselected seedlings from the same family were identified as

331     duplicates, they were both targeted for re-sampling as it was unclear which of the two was true-to-type.

332

333

334

14

335          *Pedigree records verification (GenomeStudio®, Cervus, and R)*

336     Verification of pedigree records was performed by counting the Mendelian-inconsistent errors between

337     an individual and (each of) its recorded parent(s) where genotypic data was available. These errors were

338     genotypic data inconsistent with Mendel's first law, i.e., alleles present in offspring but not present in

339     either parent. First, parent-child (PC) errors between an individual and a single parent were defined as

340     genotype calls where none of the parental alleles were present in the offspring. For example, the

341     recorded offspring might be 'BB', 'B null', or 'null null' while the recorded parent was 'AA'. In this

342     example, neither the 'B' allele nor the 'null' alleles were present in the parent. Secondly, when both

343     parents were known and confirmed, the combination of the two parents' SNP data were compared to

344     the offspring's SNP data to identify parent-parent-child (PPC) errors. PPC errors were defined as

345     genotype calls where at least one allele of the offspring was not present in any of its recorded parents.

346     For example, in the case of an 'AA' x 'AA' -> 'AB' triplet, no PC error would be observed when checking

347     each parent individually, as both parents could have contributed the 'A' allele to the offspring. However,

348     combination of the two parents would create a PPC error as neither parent could have contributed the

349     'B' allele observed in the offspring.

350

351     Three ways to count Mendelian-inconsistent errors were compared. In GenomeStudio®, a

352     'Reproducibility and Heritability' (Analysis>Reports>Reproducibility and Heritability Report>with

353     Calculating Errors) was generated to obtain the number of PC and PPC errors. Mendelian-inconsistent

354     errors were calculated in the software Cervus [51] using default parameter settings. Third, an ad hoc R-

355     script (Document S2) was used to check and identify PC and PPC relationships.

356

357     The '.gtypes' ASSIsT output file was further adjusted to the following format: the first column contained

358     an individual's 'Sample ID', the second and third columns contained the individual's 'Mother ID' and

359     'Father ID', respectively, and the subsequent columns contained the individual's genotypic data. Any

360     missing parental information was set to '-'. All alleles found in the data set were defined in the

361     'AlleleList' parameter whereas characters used for missing genotypes or missing alleles were defined in

362     the 'MissGT' and 'MissAllele' parameters respectively. After loading all functions defined in the R-script,

363     the 'CheckParAll()' function was used to identify Mendelian-inconsistent errors for individuals with at

364     least one known parent in the data set. When an individual's supposed parent was not genotyped but

365     the supposed grandparents were genotyped, the grandparents-grandchild relationship was tested with

366     the AB+AA-AA test in Excel using the template provided by van de Weg and co-workers (2018) [23].

367

368     A threshold was determined for the proportion of PC errors to confirm or reject PC relations using

369     incompletely curated marker data. PC errors were counted for a thousand pairs of two random

370     individuals in the data set that did not have a (known) PC relationship and for all pairs of individuals that

371     had a known PC relationship. A separation was observed between the resulting distributions of PC errors

372     for the two sets of individuals and a midway point between both distributions was used as threshold to

373     reject parentage of an individual. Similarly, a threshold was determined to accept or reject the

374     combination of two parents; observed PPC errors were counted for previously confirmed PPC

375     relationships and a threshold set as 110% of the highest number observed PPC errors among these

376     known relationships.

377

378     In cases of missing or erroneous parent information, efforts were made to identify the missing parent

379     and, if not possible, to identify sets of possible grandparents. Hereto, all available selected material was

380     examined (ancestors, direct parents, and breeding selections). In apple and peach, the

381     'FindPosParComb()' function of the ad hoc R-script (Document S2) was used to find PC and PPC

382     relationships. The maximum number of PC errors and PPC errors to still accept a PC relationship and PPC

16

383    relationship, respectively, were set with the 'thresholdPE' and 'thresholdPPE' parameters of the

384    'FindPosParComb()' function, respectively. In cherry, the software Cervus [51] was used to count these

385    errors and determine possible parents using the default parameter settings. When no second possible

386    parent was found in the data set, possible grandparents were identified in Excel using the template

387    provided by van de Weg and co-workers (2018) [23]. Historic records (e.g., location and time of origin) of

388    possible grandparents were checked to ensure feasibility. Furthermore, deduced grandparent-

389    grandchild relationships were only kept if they did not lead to a large number of reported errors during

390    the rest of the workflow.

391

392    Pedigree information was then updated in various input files and in GenomeStudio® (Analysis>Edit

393    Parental Relationships; then choosing individual and correct parents from drop-down menu) for further

394    analyses. All statistics in GenomeStudio® were updated when prompted.

395

396         *Genotyping error detection and adjustment (GenomeStudio®, FlexQTL$^{TM}$, and Visual FlexQTL$^{TM}$)*

397    Genotyping errors were divided in two classes: Mendelian-inconsistent errors and Mendelian-consistent

398    errors [10]. Unlike Mendelian-inconsistent errors, Mendelian-consistent errors are errors that do not

399    infringe upon Mendel's first law: a child's false allele call is present in one of the parents, but results in

400    problematic co-segregation patterns that show unexpected double recombination between markers

401    with successive genetic/physical positions. These double recombinations might be due to issues in

402    ploidy, calling, marker ordering, or phasing or, occasionally, gene conversion [10] (Document S3).

403

404    For individuals with verified pedigree relationships, remaining Mendelian-inconsistent errors were

405    detected using GenomeStudio® and FlexQTL$^{TM}$ v0.99130. In GenomeStudio®, the 'SNP Table' was filtered

406    for SNPs with Mendelian-inconsistent errors, the 'Error Table' was used to identify individuals with

17

407     Mendelian-inconsistent errors, and the 'SNP Graph' was used to examine the reported errors. FlexQTL$^{TM}$

408     input files were prepared using FlexQTL DataPrepper v1.0.0.4

409     (https://www.wur.nl/en/show/FlexQTL.htm). Three input files were needed to run FlexQTL

410     DataPrepper: a map file, a pedigree file, and a data file. The map file was obtained by adjusting the

411     ASSIsT map input file as follows: Column names were changed to 'MarkerId', Group', and 'Position' and

412     the file was saved as a comma-delimited file (.csv). The pedigree file was obtained by adjusting the

413     ASSIsT pedigree input file as follows: column names were changed to 'Name', 'Parent1', and 'Parent2'

414     and the file was saved in the '.csv' format. The data file was obtained by converting the

415     'FlexQTLDataPrepper' from ASSIsT to the '.csv' format. The data file (.dat) generated by FlexQTL

416     DataPrepper was adjusted to ensure all individuals had either both parents specified or none. Any

417     individual that had only one known parent was given a dummy parent. These dummy parents, as well as

418     any named parent not in the data set, were added to the data input file with all their genotypic data set

419     to missing. FlexQTL$^{TM}$ was used to check for Mendelian-inconsistent errors (parameter settings in Table

420     S4C). Briefly, FlexQTL$^{TM}$ was run through using an early stop ('pedimapV' parameter set to '2'; to stop

421     after checking the data for inconsistencies) and allowing for segregation distortion ('MSegDelta'

422     parameter set to 1). This analysis summarized for each marker and each individual how many

423     Mendelian-inconsistent errors were observed in the 'mconsistency.csv' file.

424

425     Mendelian-consistent errors were detected by examining double recombinations detected over small

426     regions (<10 cM) as reported by FlexQTL$^{TM}$ and Visual FlexQTL$^{TM}$. Parameter settings of FlexQTL$^{TM}$ to

427     check for double-recombinations were the same as for Mendelian-inconsistent errors above (Table S4C).

428     The FlexQTL™ output file named 'DoubleRecomb.csv' listed all singletons (single markers involved in a

429     double recombination) in the data set. Visual FlexQTL$^{TM}$ instead identifies all double recombinations

430     (including singletons) that occur within a given genetic distance. The default for this distance was 10 cM

18

431    and could be changed under 'Tools>Calculate>(Re-)Compute recombination sequences'. The report on

432    double recombinations was created through 'Tools>Export>Export recombination sequence file' which

433    provided an output file called 'DoubleRecombinations.csv'.

434

435    Genotype calls of SNPs with Mendelian-inconsistent errors or SNPs involved in detected double

436    recombinations were further examined in GenomeStudio® using the 'SNP Graph'. Where incorrect

437    cluster identification was detected, clusters were manually called using the 'SNP Graph' and FlexQTL$^{TM}$

438    was run again to ensure errors were resolved. Individuals belonging to a single cluster were chosen

439    using the 'Lasso Mode' of the 'SNP Graph'. After 'right-clicking' on the 'SNP Graph', the 'Define X Cluster

440    Using Selected Samples' was chosen where 'X' was the appropriate genotype cluster ('AA', 'AB', or 'BB').

441    The few SNPs that could not have their genotype clusters assigned simultaneously in GenomeStudio®

442    (e.g., because clusters were too closely positioned; one of the clusters for homozygous individuals was

443    between x=0.4 and x=0.6, which is true for part of the paralogous SNP one of the homozygous clusters

444    according to the ASSIsT Reference Manual p14 [34]; or because null alleles were present) were

445    genotyped as follows. Individuals belonging to a single cluster were selected using the 'Lasso Mode' of

446    the 'SNP Graph' in GenomeStudio®. 'Sample_IDs' of the chosen individuals were transferred to Excel by

447    highlighting the 'Sample_ID' column in the 'Sample Table', using the 'copy' function of the 'Samples

448    Table', and pasting them into Excel. In Excel, the copied 'SampleIDs' were then assigned a genotype call.

449    This process was repeated until all individuals had their genotype assigned. If genotype calls could not

450    be accurately made, the SNP was considered to have failed and removed from the data set.

451

452    Identification of Mendelian-inconsistent and Mendelian-consistent errors were also performed at the

453    haplotype level, conducted as described above at the single SNP level. Where an unidentified error in

454    SNP genotype scoring was detected, the corresponding SNP genotype calls were adjusted. If the calling

455     error occurred in a single or few individuals, haplotypes were manually adjusted to reflect the change in

456     SNP allele. In the rare event that a large group of individuals had their SNP genotype calls adjusted, the

457     corresponding haplotypes were re-determined using PediHaplotyper [36]. Where Mendelian-

458     inconsistent errors were due to missing SNP alleles, the individual was compared to its parent and

459     offspring to determine the correct haplotype. For example, if an individual had a SNP haplotype of 'A-?-

460     B-A' and the haplotype was not present in either parent, but a parent had a haplotype of 'A-A-B-A' and

461     no haplotype of 'A-B-B-A', the haplotype of the offspring would be set to 'A-A-B-A'. If both 'A-A-B-A' and

462     'A-B-B-A' were present in the parent, information of flanking, linked haplotypes were checked to assess

463     if the offspring's haplotype could be determined by minimizing the number of recombinations. Where

464     inconsistencies in selected material were suspected to be due to a recombination in an ungenotyped

465     progenitor, the haploblock was split in two at the suspected recombination site to avoid tracking in

466     downstream genetics analyses of recombination in selected material. The haplotypes for those two new

467     haploblocks were determined again using PediHaplotyper.

468

469     *Map error detection and adjustment (FlexQTL^TM, Visual FlexQTL^TM, and Microsoft Excel)*

470     Where double recombinations were observed and these recombinations were not due to incorrect

471     genotype scoring, a graphical genotyping approach was used to examine and possibly adjust SNP order

472     in the genetic map [52]. Graphical genotyping plots were created starting from the 'SIP_Population.csv'

473     output file of FlexQTL^TM (Document S3). FlexQTL^TM was run again to ensure the errors were resolved and

474     only if the adjustment of the SNP order did not lead to new double recombinations, a change in order

475     was accepted. SNPs were removed from the data set if they had unexpectedly high incidences of double

476     recombinations that could not be resolved by repositioning the SNPs in the map. Additionally, where a

477     SNP mapped to multiple locations in different families, the SNP was removed from the data set.

478

20

479      *Haploblock and haplotype determination (FlexQTL^{TM}, Visual FlexQTL^{TM}, and PediHaplotyper)*

480      Haploblocks were defined as regions in which no recombination was observed for selected material. For

481      phasing, parental information in the data input file of FlexQTL^{TM} was adjusted so that the pedigree was

482      trimmed to remove intermediate progenitors without genotypic data unless they were represented by

483      more than four direct offspring. Because Visual FlexQTL^{TM} does not consider any individual without

484      offspring (e.g., new breeding selections) in haploblock determination, dummy offspring with missing

485      genotypic data were added for individuals that did not have any offspring in the data set yet whose

486      recombinations were desired to contribute to determination of haploblock borders. The data was

487      phased using FlexQTL^{TM} (parameter settings in Table S4D). Next, Visual FlexQTL^{TM} was used to define

488      haploblock borders under 'Tools>Export>Export haplotype blocks file', creating the 'HaploBlocks.map'

489      file that assigns each marker to a haploblock and could be used as input for PediHaplotyper.

490

491      For SNP phasing within haploblocks, the pedigree had to be trimmed as in haploblock determination to

492      remove intermediate progenitors without genotypic data unless they were represented by more than

493      four direct offspring. However, dummy offspring introduced for haploblock determination were

494      removed again before phasing the data. FlexQTL^{TM} was then run again (parameter settings in Table S4D),

495      with the output file named 'mhaplotypes.csv', which was used as an input for PediHaplotyper.

496

497      The PediHaplotyper package [36] was loaded into R and the working directory was set to the location of

498      the input files created above ('HaploBlocks.map', 'mhaplotypes.csv', 'flexqtl.par', and 'flexqtl.sort'). In R,

499      the function 'fq_haplotyping_session(sessionID='prefix", mapfile="HaploBlocks.map")' was used to

500      create the haplotype output files in the working directory where 'prefix' was user-defined text that

501      prefixed all output file names. The 'prefix_hballeleles.dat' output file listed the composition of each

502      haplotype of each haploblock and the 'prefix_flexqtl.dat','prefix_flexqtl.map', and 'prefix_flexqtl.par'

21

503     output files were used as input files for FlexQTL™ for further data curation of the haplotyped data sets

504     (resolving both Mendelian-inconsistent and Mendelian-consistent errors as described under

505     *'Genotyping error detection and adjustment'*).

506

507         *SNP classification*

508     A SNP classifications system was established to track clustering issues and minimize future curation of

509     new data. SNPs that passed the filter criteria from ASSIsT and that were included in the final data set

510     were classified into four types: type 1 SNPs had no or less than 5% call editing during the curation

511     process and no additional genotype clusters were present; type 2 SNPs had an incorrect automated

512     cluster identification of one of the genotype clusters (e.g., 'AA' cluster called as 'AB'), showed no

513     additional clusters, and could easily be corrected; type 3 SNPs showed additional clusters because of

514     alleles with differential intensity signals but individuals could easily be called correctly; and type 4 SNPs

515     had null alleles but individuals with null alleles could be distinguished easily from true homozygous

516     individuals. Type 5 SNPs could be accurately called but their genetic or physical position could not be

517     determined accurately and were not included in the map and final data set. Type 6 SNPs were

518     monomorphic across all individuals. Type 7 SNPs were those considered as 'Failed' by ASSIsT or were

519     removed during the workflow because their genotype calls could not be manually resolved.

520

521     *Workflow creation and implementation*

522     A workflow was constructed by identifying necessary steps of data curation and ordering them in such a

523     way that the amount of time needed for data curation is minimized at each step. Thus, errors addressed

524     first were those relatively easy to identify and resolve and otherwise expected to cause problems at

525     multiple steps. The workflow was an outcome of efforts in RosBREED and FruitBreedomics on data

526     curation in apple, peach, and cherry. Statistics at each step of curation were determined from

22

527    implementing this workflow on the RosBREED germplasm described in the 'Plant Material' section

528    above.

529

530    **Results**

531

532    *Steps of the data curation workflow*

533    Initial error-detection resulted in a list of possible causes for each type of detected errors (Table 1). This

534    list identified which issues had to be resolved first and as such resulted in the workflow described below

535    (Figure 1, Document S3). The workflow developed had three main parts, each with multiple steps. The

536    first main part ensures that genetic principles can be applied, the second main part applies these

537    principles on a single marker level, and the last main part applies these principles at the haploblock

538    level. The proposed steps within each main part are described below, as conducted for apple, peach,

539    and sweet cherry.

540

541

542        Table 1: Errors observed during the curation process and their possible causes. Causes that

543    should be (mostly) already resolved by the stage a researcher would start checking for specific errors are

544    in parentheses and grey font.

| Error | Cause | Solution |
|---|---|---|
| Low call rate and impossible cluster identification | Probe binding issues | Remove SNP from data set |
| Unexpected B-allele frequencies | *(Probe binding issues)* | *(Remove SNP from data set)* |
| | Unexpected ploidy | Remove sample from data set |
| | Low sample quality | Remove sample from data set |
| High number P(P)C errors | *(Probe binding issues)* | *(Remove SNP from data set)* |
| | *(Low sample quality)* | *(Remove sample from data set)* |
| | Incorrect pedigree | Adjust pedigree record |
| | Incorrect clustering | Manually determine genotype clusters |

| | | |
|---|---|---|
| | Incorrect genotype call(s) not due to cluster issues | Adjust genotype call(s) or remove SNP from data set |
| Low number P(P)C errors | *(Probe binding issues)* | *(Remove SNP from data set)* |
| | *(Low sample quality)* | *(Remove sample from data set)* |
| | *(Incorrect pedigree)* | *(Adjust pedigree record)* |
| | Incorrect clustering | Manually determine genotype clusters |
| | Incorrect genotype call(s) not due to cluster issues | Adjust genotype call(s) |
| High number double recombinations | *(Probe binding issues)* | *(Remove SNP from data set)* |
| | *(Low sample quality)* | *(Remove sample from data set)* |
| | *(Incorrect pedigree)* | *(Adjust pedigree record)* |
| | *(Unexpected ploidy)* | *(Remove sample from data set)* |
| | Incorrect clustering | Manually determine genotype clusters |
| | Incorrect marker position in map | Adjust marker position or remove marker if it cannot be accurately mapped |
| | Incorrect genotype call(s) not due to cluster issues | Adjust genotype call(s) |
| | Incorrect phasing | Find responsible individual and make genotype missing |
| Low number double recombinations | *(Probe binding issues)* | *(Remove SNP from data set)* |
| | *(Low sample quality)* | *(Remove sample from data set)* |
| | *(Incorrect pedigree)* | *(Adjust pedigree record)* |
| | *(Incorrect clustering)* | *(Manually determine genotype clusters)* |
| | *Nearby double recombination\** | Resolve nearby double recombination |
| | Incorrect marker position in map | Adjust marker position or remove marker if it cannot be accurately mapped |
| | Incorrect genotype call(s) not due to cluster issues | Adjust genotype call(s) |
| | Incorrect phasing | Wait for haploblock analysis to resolve issue |
| Incorrect haplotype determination | *(Probe binding issues)* | *(Remove SNP from data set)* |
| | *(Low sample quality)* | *(Remove sample from data set)* |
| | *(Incorrect pedigree)* | *(Adjust pedigree record)* |
| | *(Incorrect clustering)* | *(Manually determine genotype clusters)* |
| | *(Incorrect marker position in map)* | *(Adjust marker position or remove marker if it cannot be accurately mapped)* |
| | *(Incorrect genotype call(s) not due to cluster issues)* | *(Adjust genotype call(s))* |
| | Incorrect phasing | Manually correct phasing (determine correct haplotypes) |
| | Recombination within haplotype | Adjust haploblock borders |

*Nearby double recombination can occur for two adjacent markers with many double recombinations and markers with few double recombinations. However, nearby double recombinations rarely lead to a high number of double recombinations for a single marker

545

546

547     Figure 1: Steps of the high-resolution genotypic data curation workflow to ensure a quick and

548     efficient curation process. Steps that identify errors are shown in white boxes; procedures needed for

549     detecting, keeping track of, and resolving errors but do not identify errors directly are in grey boxes.

550     After obtaining a first set of genotypic data, initial steps ensure that inheritance principles can be readily

551     applied by removing individuals and markers that do not follow these principles and by ensuring

552     pedigree records are correct. In the next set of steps, inheritance principles are applied at the individual

553     marker level. In the final set of steps, these principles are applied at the haploblock level. Output used to

554     detect and resolve observed errors at each step are given in italics. The leaf symbol indicates errors at

555     the level of individual; the intensity plots symbol indicates errors at the level of SNP scoring; the genetic

556     map symbol indicates errors at the level of genetically linked markers and phased alleles. When applying

557     inheritance principles in parts 2 and 3, alleles that do not occur in an individual's parents ('Mendelian-

558     inconsistent errors') are first resolved before addressing remaining genotyping errors ('Mendelian-

559     consistent errors'). Several procedures, such as marker call adjustments and map order adjustments, are

560     performed throughout the steps of the workflow to resolve errors detected. Each time after performing

561     these common procedures, specific steps of the workflow must be repeated, forming an iterative

562     process that ends when all errors are resolved.

563

564

565

566

567    *1. Ensuring inheritance principles can be applied*

568    After creating an initial data set of genotypic data set in GenomeStudio®, a first set of analyses was

569    performed. Because genotypic errors are identified based on principles of inheritance in diploids,

570    individuals and markers that do not to follow these principles had to be removed first (Figure 1). When

571    doing so, individuals with unexpected intensity patterns had to be removed first (Figure 1) as they were

572    influencing the clustering of all individuals in the germplasm. Individuals with poor quality DNA were

573    usually poorly genotyped, resulting in many data inconsistencies. Additionally, polyploids (individuals

574    having one or more additional full chromosome sets) and aneuploids (individuals having an irregular

575    number of copies for one or more chromosomes) were expected to have intensity ratios for

576    heterozygous loci that differed from diploid individuals. Removal of individuals with poor DNA quality

577    and suspected polyploids and aneuploids was observed to improve genotype cluster definitions and

578    thereby the genotype calling of remaining individuals.

579

580    Once individuals with ploidy and sample quality issues were removed, a set of well performing markers

581    had to be obtained (Figure 1). Markers with unreliable scoring were observed to lead to many

582    inconsistencies in subsequent steps. Thus, their early removal would ensure that a relatively low

583    number of inconsistencies remained in the data set, expected to greatly reduce the observed

584    inconsistencies and time needed for further steps.

585

586    Identifying and correcting incorrect PC and PPC relationships was a prerequisite to using pedigree

587    information for the identification of marker calling errors in each data set. Imposing principles of

588    inheritance on actually unrelated individuals led to many false errors at the marker and map level.

589    Conversely, identifying thus far unknown PC and PPC relations helped to identify errors at the marker

590    and map level elsewhere in the data set and was expected to improve the power of downstream QTL

591    analyses. Thus, recorded pedigree information needed to be validated and previously unknown pedigree

592    relationships deduced before curating individual marker calls and marker order errors (Figure 1).

593    Duplicate individuals were also detected at this stage as they could help resolve sampling errors and

594    reduce the number of individuals needing detailed error-checking.

595

596

597            *2. Applying inheritance principles at the marker level*

598    When Mendelian-inconsistent errors were present, at least one allele was incorrect. This issue had to be

599    resolved before the (corrected) allele could be phased with the alleles of flanking markers. Otherwise,

600    even the other allele, which might have been correct, could have been incorrectly phased with the

601    alleles of flanking markers, causing additional observed but false recombinations. Thus, to minimize the

602    time required to resolve Mendelian-consistent errors by investigating many supposed double

603    recombinations, Mendelian-inconsistent errors had to be addressed first.

604

605    Markers with a high number of errors were investigated before markers with a relatively low number of

606    errors among progenitors. Then, markers with a low number of errors for seedlings were investigated as

607    they were expected to have the least effect on the remaining data set.

608

609    Any supposed double recombinations that occurred at the same region in multiple individuals had to be

610    resolved first as they were very unlikely, could be due to a single error, and could influence a large set of

611    individuals. Next, suspicious double recombinations that occurred over multiple loci in ancestors had to

612    be checked, followed by singletons in ancestors. Finally, singletons in seedlings were checked, but they

613    were expected to be the least harmful when incorrect because of little to no effect on the remaining

614    data set.

27

615

616     When no genotype calling or map errors were detected, phasing errors were investigated by checking

617     the phasing of individuals that shared the parent whose homolog was observed to have a double

618     recombination. In the rare case that incorrect phasing by FlexQTL™ led to a double recombination in

619     multiple individuals of a single family or parent, it was always caused by one or two individuals in which

620     the position of (a single) recombination was incorrectly determined. In those cases, individual(s) for

621     which the SNP was involved in a single recombination had their genotype set to missing. This adjustment

622     led to correct phasing of all other individuals and removal of reported double recombinations. Double

623     recombinations that were observed in a single individual and that were not due to incorrect genotype

624     clustering or incorrect map positions were accepted as the result of true double recombination events.

625

626            *3. Applying inheritance principles at the haploblock level*

627     Haploblock and haplotype determination was based on correctly identifying recombinations through

628     correct phasing across generations and combining individual SNP alleles into haplotypes. Thus, any

629     remaining errors at the SNP level or map level were expected to lead to errors in haploblock and

630     haplotype determination. Therefore, all observed inconsistencies at the individual SNP level had to be

631     resolved before inconsistencies were detected at the haploblock level. The genetic principles applied

632     throughout the workflow are expected to also hold up at the haploblock level and therefore haplotypes

633     had to be checked for Mendelian-consistent errors and Mendelian-inconsistent errors.

634

635

636

637

638

639     *Implementation of the workflow on RosBREED apple, peach, and sweet cherry germplasm*

640

641          *1a. Removing samples: non-diploid individuals and low-quality samples*

642     In apple, the 'B allele frequency' plot of 744 of the diploid individuals (80.7 %) was very close to that

643     expected for diploid individuals (Figure 2A; Table S1) and results of these diploid individuals were

644     considered to be of good quality. Another 71 individuals (7.7%) showed some variation from the

645     expected B allele frequency, especially for homozygous SNPs, but the three genotypes could be easily

646     distinguished (Figure 2B; Table S1) and their results quality was considered to be intermediate. Finally,

647     107 (11.6%) had 'B allele frequency' plots that showed a wide variation around the expected frequency

648     (Figure 2C; Table S1) and their results quality was considered to be bad. No individuals with bad quality

649     results were found for peach or sweet cherry.

650

651

652          Figure 2: Histograms of B-allele frequency (left) and B-allele frequency for each SNP plotted

653     against its genomic position (right). Such histograms were used to assess a sample's genotyping quality

654     and ploidy. Examples shown are of a sample with good quality genotype calls (panel A), with

655     intermediate quality of genotype calls (B), with bad quality of genotype calls (C), and that is triploid (D).

656

657

658     For apple, most individuals with poor quality results had their DNA extracts transported outside the U.S.

659     for genotyping and the poor results were suspected to be caused by a reduction in DNA quality due to

660     the delay in clearing customs, while only nine individuals with poor quality were from those genotyped

661     in the U.S. The call rate in GenomeStudio® differed between the individuals that had good, intermediate,

662     or bad quality, with the call rate dropping as the level of quality lowered (Figure S2).

29

663

664    For apple, five triploid individuals were identified (Table S1). One was the known triploid cultivar

665    'Jonagold' while the others were unselected seedlings (Table S1; Figure S1A). Two other unselected

666    seedlings had their B-allele frequencies divided over 5 clusters of the GenomeStudio® plot, which

667    indicated they could be tetraploid or a mixture of two samples (Table S1; Figure S1B). No aneuploids

668    were detected in the apple germplasm. However, one individual from the Crop Reference Set, 'AE213-

669    200' and one individual of a Breeding Pedigree Set were identified as segmental aneuploids (missing one

670    copy of a large chromosomal segment). They were undetectable in the B-allele frequency analysis and

671    instead identified by a relatively large number of PC errors and double recombinations observed for only

672    that chromosomal segment. No polyploids, aneuploids, or segmental aneuploids were detected in peach

673    and sweet cherry.

674

675    The final number of individuals used in the rest of the workflow was 835, 621, and 528 for apple, peach,

676    and sweet cherry, respectively, consisting of 139, 48, and 56 direct parents of full-sib families, ancestors,

677    and cultivars, 76, 24, and 9 selections and 620, 548, and 463 unselected seedlings over 45, 26, and 41

678    families of 4–62 full-sibs, respectively (Tables S1-S3).

679

680        *1b. Obtaining a set of reliable SNPs*

681

682    A subset of SNPs with reliable genotyping scores was obtained using ASSIsT (Table 2). Although

683    discarded by ASSIsT, SNPs from the 'Two cluster SNPs' category were retained as many of them were

684    considered to contain useful information. A total of 4636 (59%), 6098 (75%), and 1727 (30%) of the SNPs

685    on the apple, peach, and cherry arrays, respectively, were maintained after filtering. Subsequent steps

686    of the workflow reduced the number of SNPs in the final data set further to 3855, 4005, and 1617 for

687 apple, peach, and sweet cherry, respectively. Thus 83%, 66%, and 91% of the SNPs retained after using

688 ASSIsT for apple, peach, and sweet cherry, respectively, resulted in high-quality data.

689

690

691 Table 2: Summary of SNP classification by ASSIsT for apple, peach, and sweet cherry. SNP

692 classifications are grouped in retained and discarded SNPs.

| SNP classification | Apple | Peach | Sweet Cherry |
|---|---|---|---|
| **Retained SNPs** | | | |
| *Robust SNPs* | | | |
| Robust | 1435 | 743 | 373 |
| OneHomozygRare_HWE | 368 | 62 | 109 |
| OneHomozyRare_NotHWE | 369 | 188 | 161 |
| DistortedAndUnexSegreg | 1364 | 3696 | 555 |
| *Other* | | | |
| Two cluster SNPs | 1100 | 1409 | 529 |
| *Total* | *4636* | *6098* | *1727* |
| **Discarded SNPs** | | | |
| NullAllele-Failed | 57 | 145 | 43 |
| Monomorphic | 1307 | 1057 | 3478 |
| Failed | 2888 | 844 | 448 |
| *Total* | *4252* | *2056* | 3969 |
| **Total** | **8888** | **8144** | **5696** |

693

694

695 *1c. Correcting pedigree information and identifying duplicates*

696 The number of PC errors in apple between two randomly paired individuals without PC relationship

697 averaged 195, with a minimum of 17 (comparison between two full-sibs) and 99% of these comparisons

698 had more than 40 errors. In contrast, average and maximum number of PC errors between two related

699 individuals with a known PC relationship was 2 and 17, respectively, and 99% of these comparisons had

700 less than 10 PC errors. The threshold to reject a PC relationship was set at 23 errors, which roughly

701 corresponded to 0.5% of total markers. For 103, 66, and 22 individuals, one recorded parent was

702    incorrect in apple, peach, and sweet cherry respectively, and for 36, 14, and zero individuals, both

703    recorded parents were incorrect. For 106, 1, and 19 of these individuals in apple, peach, and sweet

704    cherry, one or both of the true parent(s) was found within the germplasm set. The final number of

705    generations spanned by the corrected pedigrees was eight, nine, and six for apple, peach and sweet

706    cherry, respectively.

707

708        *2a. Finding Mendelian-inconsistent errors at the SNP level*

709    FlexQTL™ summarized the number of Mendelian-inconsistent errors for each marker and each

710    individual. In GenomeStudio®, the 'SNP Table' would summarize the number of P(P)C errors for each

711    SNP and a separate 'Error Table' had to be consulted to determine which individuals were involved in

712    these errors. FlexQTL™ mostly reported the error under the parent, the R-script reported the error

713    under the offspring, and the 'Error Table' of GenomeStudio® reported the genotypes of both parent(s)

714    and offspring. As a consequence, errors between a single parent and multiple of its offspring would be

715    reported as one erroneous (parental) genotype in FlexQTL™ whereas GenomeStudio® reported the

716    error for each offspring. However, FlexQTL™ did identify errors between grandparents and

717    grandchildren when the missing parental genotype could be imputed.

718

719    FlexQTL™ detected 1209, 2230, and 686 Mendelian-inconsistent errors distributed over 541, 760, and

720    42 SNPs in apple, peach, and sweet cherry respectively. In apple, GenomeStudio® detected 10,201 PC

721    errors and PPC errors over 2303 SNPs. Although GenomeStudio® identified which pairs of individuals led

722    to these errors, some of the detected Mendelian-inconsistent errors did not occur in the data set due to

723    differences in genotype scoring between ASSIsT and GenomeStudio®. Before removal of these

724    Mendelian-inconsistent errors, 41,717, 29,009, and 2505 double recombinations involving a single

725    marker were detected in FlexQTL™ in apple, peach, and sweet cherry, respectively, through the

32

726     'DoubleRecomb.csv' file, whereas only 6177, 4905, and 1739, respectively, of these recombinations

727     were observed after removal of all Mendelian-inconsistent errors.

728

729     *2b. Identifying Mendelian-consistent errors at the SNP level*

730     Most double recombinations that occurred in the same genomic region in many individuals could be

731     resolved by adjusting incorrect marker calls. A total of 648, zero, and 209 markers in apple, peach, and

732     sweet cherry, respectively, had one or more of their genotype calls adjusted to resolve double

733     recombinations. Most other double recombinations that occurred in multiple families could be resolved

734     by repositioning the marker in the genetic map using a graphical genotyping approach. In total, 115,

735     zero, and zero ### SNPs were moved from their original position in the map to resolve double

736     recombinations for apple, peach, and sweet cherry, respectively. Many recombination events that

737     occurred in a single or few individuals over a single marker were resolved by first resolving the double

738     recombinations that occurred in many individuals. Most of the remaining double recombinations were

739     solved by either changing single incorrect genotype call or adjusting marker order in the map. Only a few

740     phasing issues were observed where (almost) all offspring of a founder showed a double recombination

741     that could be resolved by adjusting the phase of the alleles in that founder. A total of 15, 156, and 63

742     markers were discarded for apple, peach, and sweet cherry, respectively, because they led to

743     unresolvable map issues. The total number of remaining reported singletons was 68, 47, 51 for apple,

744     peach, and sweet cherry, respectively, and these were considered to be true double recombinations.

745

746     During data curation, genetic maps were generated for each crop (Tables S5-S7) by adding new SNPs to

747     existing maps, by converting physical positions into genetic positions, and/or by updating initial genetic

748     positions to minimize the number of double recombinations. For apple, 885 SNPs were added and 658

749     previously-mapped SNPs were removed as they did not perform well in our wider germplasm. Addition

33

750     of SNPs at the chromosome ends enlarged the original map by 7 cM. The resulting apple map was 1179

751     cM long with chromosome lengths ranging from 57.6 cM (linkage group (LG) 6) to 103.6 cM (LG 15). The

752     number of SNPs on each LG ranged from 167 SNPs on LG 6 to 359 SNPs on LG 2. The genetic map of

753     peach was 893.2 cM long; LG 5 was the shortest (72.9 cM) and LG 1 was the longest (190.2 cM). The

754     number of SNPs on each LG ranged from 294 on LG 5 to 772 on LG 4. In sweet cherry, chromosome

755     lengths ranged from 56.8 cM (LG 5) to 141.2 cM (LG 1), with a total map length of 655.4 cM. The

756     number of SNPs on each LG ranged from 137 on LG 5 to 350 on LG 1.

757

758          *3. Determining and resolving errors for haploblocks and haplotypes*

759     The genetic maps of apple, peach, and sweet cherry were at first divided in 840, 103, 132 haploblocks,

760     respectively, within which no recombination was observed in selected germplasm. After haplotype

761     generation, 1262, 2012, and 74 Mendelian-inconsistent errors were reported by the mconsistency.csv

762     file generated by FlexQTL$^{TM}$. An additional 124, 429, and 64 recombinations were detected within the

763     haploblocks for selected germplasm, resulting in the generation of additional haploblocks. The

764     remaining Mendelian-inconsistent errors were mostly due to missing data within a haplotype that could

765     not be resolved automatically. This missing data within haplotypes led to the assignment of haplotype

766     numbers that were different to parental haplotypes that were therefore perceived as errors. In addition,

767     some inconsistencies between SNP data and haplotype data were observed after haplotype generation

768     that were easily resolved by looking at the 'SNP Graph' in GenomeStudio® and adjusting either the

769     haplotype or the SNP call.

770

771     The final number of haploblocks was 964, 135, and 196 for apple, peach, and sweet cherry respectively.

772     For apple, the genetic length of the haploblocks varied between 0 and 7.77 cM with an average of 0.3

773     cM, the haploblocks contained between 1 and 15 SNPs, and the haploblocks contained an average of 4

34

774    SNPs. The number of haploblocks per apple LG ranged from 42 on LG 6 to 79 on LG 15, with an average

775    of 57 haploblocks per LG. In peach, the length of the haploblocks varied between 0 cM and 30.47 cM

776    with an average of 5.8 cM, the haploblocks contained between 1 and 210 SNPs, and the haploblocks

777    contained an average of 30 SNPs. The number of haploblocks per peach LG ranged from 7 on LG 5 to 37

778    on LG 4, with an average of 17 haploblocks per LG. For sweet cherry, haploblocks had an average length

779    of 2.6 cM, with a minimum of 0 cM and a maximum of 15.0 cM. The average number of SNPs per sweet

780    cherry haploblock was 8, with a minimum of 1 and a maximum of 61 SNPs. The average number of

781    haploblocks per sweet cherry LG was 24, with a minimum of 16 haploblocks on LG 5 and LG 7 and a

782    maximum of 47 haploblocks on LG 1.

783

784

785         *SNP classification system*

786    The final number of SNPs in the haplotyped data set was 3858, 4005, and 1617 for apple, peach, and

787    sweet cherry, respectively. A total of 3350 (87%), 4005 (100%), and 1610 (99.6%) of these SNPs were

788    classified as type 1 SNPs, which ultimately needed editing for less than 5% of their genotype calls in

789    apple, peach, and sweet cherry, respectively (Tables S8-10). Type 2 SNPs, for which genotype clusters

790    were shifted, totaled 300 (8%), zero, and seven (0.4%) SNPs for apple, peach, and sweet cherry,

791    respectively, and this shift in cluster position lead to incorrect identification of one of the three clusters

792    in the original automatic clustering by GenomeStudio®. Type 3, SNPs with additional clusters, were

793    assigned to 80 (2%), zero, and zero SNPs in apple, peach, and sweet cherry, respectively, and this

794    presence of additional clusters led to incorrect genotype scoring of these SNPs that required subsequent

795    curation. Type 4, SNPs with null alleles, were assigned to for 125 (3%), 145 (excluded from the final data

796    set), and 43 (excluded from the final data set) SNPs in apple, peach, and sweet cherry, respectively, and

797    these null alleles prevented correct automatic scoring for some individuals.

35

798

799    **Discussion**

800    We established a workflow to efficiently and confidently identify and remove genotyping errors from

801    genotyped and pedigreed germplasm sets for apple, peach, and sweet cherry. The proposed workflow

802    (Figure 1) enables directed identification of markers and individuals with genotyping errors. It uses

803    simple genetic principles such as inheritance of parental alleles, the co-segregation of linked markers,

804    and the likelihood of double recombinations to find these errors. The order of steps was determined to

805    efficiently minimize errors found in later steps and thereby minimize overall time needed to find errors

806    in the data set. For example, in apple, any incorrect PC relationship would lead to an average of 196

807    reported Mendelian-inconsistent errors, and any unresolved Mendelian-inconsistent errors led to an

808    average of 30 more reported Mendelian-consistent errors. The developed workflow was demonstrated

809    on Illumina SNP array data and some software is specific to this platform, but the same workflow order

810    and genetic principles are appropriate for other marker types and genotyping platforms. The workflow is

811    especially useful when medium- and high-throughput genotyping tools are used for which checking each

812    individual marker would be too time-consuming.

813

814

815        Table 3: Recommended software for each step of the genetic marker data curation workflow

816    when using Illumina Infinium® SNP arrays.

| Workflow step | Recommended software |
| --- | --- |
| Identify polyploids, aneuploids, and samples with low quality | GenomeStudio® to obtain B-allele frequencies, R to plot B-allele frequency for each sample |
| Create subset of reliable SNPs | ASSIsT |
| Identify duplicate samples | PLINK |
| Identify incorrect P(P)C relationships | GenomeStudio® |
| Identify unknown P(P)C relationships | R |
| Identify unknown grandparent-grandchild relationships | Excel* |

| | |
|---|---|
| Identify and resolve (remaining) Mendelian-inconsistent errors | GenomeStudio®, FlexQTL™ |
| Identify and resolve Mendelian-consistent errors | Visual FlexQTL™ + GenomeStudio® |
| Identify and correct map order inconsistencies | Visual FlexQTL™ |
| Identify phasing issues | FlexQTL™ + Visual FlexQTL™ |
| Haploblock border determination | Visual FlexQTL™ |
| Haplotype determination | |
|    - Phasing | FlexQTL™ |
|    - Haplotype assignment | PediHaplotyper |
|    - Curation (automated) | FlexQTL™ |

817    * Template in Suppl. File 1 of Van de Weg and co-workers (2018) [23]

818

819

820    *Order and considerations of workflow steps*

821    Different types of errors can be present in genotypic and pedigree data, caused by different kinds of

822    issues (Table 1). To minimize the time needed for curation of these data, the proposed error checks

823    need to be performed in a specific order. By first tackling issues that are common for many types of

824    errors, subsequent curation of remaining errors becomes easier and quicker.

825

826    *Removing individuals with low quality or irregular number of chromosome sets*

827    The B-allele frequency plots provided a quick and easy way to identify and remove individuals with an

828    irregular number of chromosome sets (polyploids and aneuploids) and individuals with low DNA quality.

829    Removal of such individuals improved SNP calling and thus reduced the number of errors to be dealt

830    with in later steps. A couple of individuals with poor quality that were originally kept, because of their

831    importance as breeding parents, resulted in many PC errors. Making all their original SNP calls missing

832    enabled automated imputation of most of these data points based on genetic information of relatives.

833    Subsequent re-genotyping of these individuals matched the imputed data completely, confirming that

834    the errors observed were due to low-quality DNA samples and not to incorrect PC relationships.

835    Polyploid and aneuploid individuals did not show a higher number of P(P)C errors, as expected. In

836     contrast, these chromosome number abnormalities led to higher rates of false double recombination,

837     either genome-wide (polyploid) or local [(segmental) aneuploids], that cannot be readily resolved other

838     than by removal of these specific individuals.

839

840     The histogram function in GenomeStudio® enabled quick identification of polyploids and individuals

841     with very poor DNA samples without the need for additional steps in Excel, R, or other software.

842     However, identification of aneuploids and individuals with potentially low-quality DNA samples was not

843     as straightforward. Plotting the B-allele frequency against physical or genetic marker order (when

844     available) required additional data manipulation and generation of the plots in software outside

845     GenomeStudio®, but most of it could be automated using R and custom scripts. Therefore, we suggest

846     using GenomeStudio® for initial removal of poor-quality samples and polyploids, and afterwards, when

847     positional information for the markers is available, screening for aneuploids with the method described

848     by Chagné at al. (2015).

849

850     *Obtaining a set of reliable SNPs*

851     SNPs with major scoring issues that cannot be easily resolved manually need to be removed from the

852     data set. The early detection and removal of these unreliable SNPs greatly reduces the number of

853     marker and map errors reported, as well as the time spent evaluating these SNPs in later workflow

854     stages. By using ASSIsT, a quick subset of SNPs with robust genotype calls could be generated. On

855     average across the three crops, 80% of this subset was retained in the final data set, which is lower than

856     the 99% for single full-sib families that was reported by Di Guardo and co-workers (2015) [35]. As the

857     number of generations and full-sib families in the germplasm increase, more SNPs with null alleles are

858     likely to be detected and the more complicated the genotype calling of these SNPs can become. In turn,

38

859     this can lead to an increased discarding of SNPs, which could explain the lower proportion of SNPs

860     retained in our germplasm sets compared to that reported by Di Guardo and co-workers (2015) [35].

861

862     Markers with null alleles identified by ASSIsT were removed from the data set, as they could only be

863     identified and automatically called in specific $F_1$ families rather than in all families and across

864     generations. However, many SNPs with null alleles that were later identified in the workflow could be

865     accurately genotyped manually as long as homozygous 'AA' and 'BB' individuals could be distinguished

866     from individuals that carried a null allele. This distinction was time-consuming and therefore we

867     recommend saving these SNPs only when it justifies the time needed to do so. Examples when such

868     markers can be of high value are in the construction of genetic linkage maps, even if multiple mapping

869     populations are used [45], when they occur in a region of low coverage, or when they occur in a region

870     of specific interest and help define additional alleles.

871

872     Very few other options exist to create a subset of high-quality genome-wide markers across pedigreed

873     germplasm. GenomeStudio® does provide several quality scores that have been used before in SNP

874     filtering, but no guidelines exist on what threshold values to use. Using parameter thresholds regularly

875     reported in literature [26,53–56] (GenTrain Score > 0.7, 50GC Score > 0.4, ClusterSep Score > 0.25, Call

876     Rate > 0.9, and Minor Freq > 0.01) on the current data, the proportion of retained, unreliable, or

877     monomorphic SNPs would be 12.3%, 23.1%, and 6.7% in apple, peach, and sweet cherry, respectively,

878     and a large proportion of good SNPs would be discarded (27.8%, 28.2%, and 7.6%, respectively). Thus,

879     ASSIsT greatly increased the number of reliable SNPs that were retained without reducing the quality of

880     the subset of SNPs, making it the most efficient method to choose SNPs without prior knowledge on SNP

881     performance.

882

883 *Updating pedigree records*

884 As thresholds to confirm or discard historic pedigree information depends on the germplasm,

885 genotyping platform, and data quality, they need to be assessed case-wise. A custom R-script provided

886 quick and easy determination of the number of PC and PPC errors. However, the custom code required a

887 significant amount of time to identify possible parents when one or both parents were unknown,

888 especially for larger data sets. Similar issues were observed for Cervus, which took a long time to run

889 (days) and did identify some incorrect relationships, especially for inbred material. Cervus also requires a

890 specific data format and we experienced some problems running the software for large data sets that

891 were not immediately resolved. GenomeStudio® provided the quickest way to determine the number of

892 PC and PPC errors, which could be determined immediately after loading the raw intensity data.

893 However, new PC relationships could not automatically be determined and only SNPs retained by ASSIsT

894 should be used when using GenomeStudio® to determine the number of PC errors, to avoid inflating the

895 number of PC errors. Therefore, we recommend using GenomeStudio® to confirm existing pedigree

896 records when using Illumina arrays and using an R-script to determine new, previously unknown, PC

897 relationships. Time-consuming analyses in R could be resolved by using a subset of markers equally

898 spread across the genome. For confirming and identifying possible grandparent-grandchild relationships,

899 we recommend the Excel template provided by van de Weg and co-workers (2018) [23]. However, this

900 method can misconstrue aunts-uncles/nephew-nieces and individuals with other close relationships to

901 the target individual as grandparents. Therefore, we recommend to only use this strategy when the user

902 has a good understanding of the germplasm such as the origin of the material and the degree of

903 inbreeding.

904

905 Individuals with only one parent known can still be used in a pedigree-based approach to find errors in

906 the data set, although some errors might remain unnoticed. We recommend using the 'M_' and 'F_'

40

907    prefixes to the individual's name to designate the unknown mother or father, respectively. When it is

908    unclear whether the unknown individual is the mother or the father, the 'UP_' prefix can be used. Using

909    this system instead of a non-descriptive name such as 'dummy 1' creates a clear connection between

910    the individual with an unknown parent and the placeholder individual that is introduced. When the

911    correct parent is later found, it also allows the quick replacement of the placeholder by the correct

912    name (and corresponding genotypic data). Use of the same name for any missing parent should be

913    avoided (e.g., using 'dummy' for all missing parents) unless the missing parent is unequivocally the

914    parent of multiple individuals. If the same name is used incorrectly for multiple missing parents, the

915    genotype of that single missing parent is expected by FlecQTL™ to be consistent with inheritance

916    principles for all of its assigned offspring, potentially creating a large number of errors in further steps.

917

918    Although non-diploid individuals should be removed from the workflow before identifying reliable SNPs,

919    they can have their pedigree checked if needed. Regardless of their ploidy, individuals should only

920    contain alleles that are present in their parents. For example, a triploid individual with a marker call at

921    one SNP of 'AAA' will be scored as 'AA', but can still not have a 'BB' parent. However, caution is advised

922    as the grandparents through the parent that provided the unreduced gamete will also share a full allele

923    set with any polyploid individual and thus these grandparents could also be incorrectly assigned as a

924    parent of the polyploid individual. For example, the triploid 'Zonga' and its (diploid) grandparent 'Cox's

925    Orange Pippin' share a full allele set (through an unreduced gamete of 'Alkmene') and thus no PC errors

926    are reported [57]. However, only the combination of 'Delcorf' and 'Alkmene' could explain the

927    genotypes of the triploid 'Zhonga' (AB+AA-AA test [23]). Thus, for triploids, not only do parents and

928    offspring lead to no PC errors but some grandparents do as well, and the second parent is needed to

929    identify the true PC relationship.

930

41

931 *Creating or extending genetic maps.*

932 This study used available genetic maps for apple and cherry (i.e., [20,21,45,47]), integrated them when

933 needed, and used available physical information (from [44] and [46]) to add any markers that were not

934 already mapped. Some of these added markers were positioned at chromosome ends, which resulted in

935 the increase of the map size by 7 cM for apple. In addition, the orientation of apple chromosome 5 was

936 inverted here to match the orientation of the latest genome version [44]. If no genetic map is available,

937 one will need to be constructed alongside genotypic data curation. The need for a precise genetic

938 position of markers on the 9K peach array prompted development of consensus linkage map for peach

939 [58] that in the future could serve as a reference map to estimate genetic positions of unmapped

940 markers. A mapping approach for pedigreed, multi-parental maps is described by Di Pierro and co-

941 workers (2016) [45].

942

943 *Resolving remaining Mendelian-inconsistent errors*

944 Use of GenomeStudio® for detecting Mendelian-inconsistent errors is limited to Illumina array SNPs and

945 cannot be used for other markers or haplotypes created in later steps of the workflow. In addition, some

946 SNPs had their SNP scoring improved with ASSIsT and manual curation, and thus the genotype scoring of

947 GenomeStudio® might not reflect the actual data. Although this latter limitation is also true when

948 confirming pedigree data, the few differences in genotype calls between GenomeStudio® and ASSIsT are

949 not expected to alter the outcome of pedigree confirmation. In contrast, when resolving single

950 Mendelian-inconsistent errors, it is important to know that the error is indeed present in the data set.

951 Although Cervus counts the number of Mendelian-inconsistent errors, it does not report which markers

952 are causing issues for which individuals, making it impractical to use to remove the remaining PC and

953 PPCerrors. In contrast to GenomeStudio®, FlexQTL™ can handle multiple allele formats and is thus

954 suited for the curation of both SNP data and haplotype data. In addition, FlexQTL™ checks for

42

955     consistency over multiple generations, which enables detection of errors even if a genotype is missing in

956     an intermediate individual. It also imputes missing data whenever possible. A disadvantage of FlexQTL™

957     is that it only reports one of the two individuals, often the parent, for which an error occurred; it is then

958     up to the user to find the second individual, often the offspring, involved in the Mendelian-inconsistent

959     error. Therefore, we recommend using FlexQTL™ to identify Mendelian-inconsistent errors and

960     resolving them with the help of GenomeStudio®.

961

962     *Using map and phasing information to detect Mendelian-consistent errors*

963     FlexQTL™ performed very accurate phasing and only a few phasing issues were noticed. Most of these

964     phasing issues were observed as double recombinations in offspring of an individual that served as a

965     founder. The lack of parental info for this founder provided FlexQTL™ more freedom to phase alleles, as

966     the phasing in the founder did not need to match its parents. Incorrect phasing was most likely caused

967     by one or very few offspring for which a true recombination occurred in the map region. In those

968     individuals, no double recombination occurred, and the incorrect phasing inferred by FlexQTL™

969     minimized the interval over which the true recombination occurred. However, this minimalization of the

970     recombination interval incorrectly specified where the recombination had occurred, causing incorrect

971     phasing and resulting in one or multiple false double recombinations in full- and half-sibs of the

972     individual(s) with the true recombination. Making genotype calls missing for the individual(s) with a

973     recombination in that area enlarged the recombination interval for those individuals, but also led to

974     correct phasing in their parent and resolved the supposed double recombinations in their full- and half-

975     sibs. Very few other phasing issues were observed that could not be resolved on a single SNP level but

976     were later resolved at the haploblock level. Thus, a small number of phasing issues can be accepted

977     when moving forward to generating haploblocks and they could be nullified by FlexQTL™ by setting the

978     parameter 'DeleteDR' to 1.

979

980 *Haploblock and haplotype determination*

981 Visual FlexQTL™ showed good accuracy (between 12% and 33% of the initial haploblocks had to be

982 divided into additional haploblocks to avoid recombination within haploblocks for selected material) in

983 determining haploblock borders based on historic recombination events. Two reasons exist for not

984 identifying all historic recombinations for haploblock border determination. First, Visual FlexQTL™

985 determines the border as the middle of the recombination interval. The more non-informative markers

986 present in the recombination interval (due to homozygosity or lack of co-segregation (phase)

987 information), the less likely that the middle position is the true position of the historic recombination

988 (which determines the haploblock border). Secondly, FlexQTL™ determines haploblock borders

989 sequentially, starting with small recombination intervals; if multiple recombinations occur in the same

990 region, one haploblock border could suffice to account for all recombinations. This approach thus

991 minimizes the number of recombination sites needed to explain observed segregation data. In reality,

992 the recombinations could have occurred between different markers, requiring that region to be split in

993 additional haploblocks to avoid recombination within haploblocks for selected material.

994

995 PediHaplotyper's haplotypes did not always match with SNP data. In most cases, these inconsistencies

996 were introduced during the marker consistency check with FlexQTL™ to ensure the haplotypes in an

997 individual matched those of its parents and offspring. When the haplotype that caused the inconsistency

998 was represented well in the pedigree, the haplotype was correct and the original genotype call for the

999 SNP was incorrect. Thus, in these cases, haplotype curation identified additional errors in the SNP data.

1000 These errors were mostly caused by (very) small incorrectly identified genotype clusters or by single

1001 calling errors in the data set that were not detected earlier. When haplotypes in poorly represented

1002 individuals (one or two directly related individuals in the data set) showed an inconsistency with the SNP

44

1003    data, the SNP data was mostly correct and an error had occurred during haplotyping. The error could

1004    span multiple generations leading to inconsistencies for multiple individuals but its impact on the

1005    dataset was small as the overall representation of the incorrect haplotype was small. In the rare case

1006    that a poor representation led to incorrect haplotype determination, the actual cause of the

1007    inconsistency often remained unclear, but for some it was due to a recombination within a haploblock

1008    for an un-genotyped ancestor or one of the direct parents of such an ancestor.

1009

1010    Haploblock borders are not fixed and can change based on the application of the final data set and the

1011    germplasm used. For example, for QTL analyses some of the haploblocks defined here will be too large

1012    as they span multiple cM; they will show within-haploblock recombination in numerous unselected

1013    offspring thereby increasing the number of missing haplotype calls thus increasing uncertainty in QTL

1014    position (including the widening of QTL intervals). Haploblock sizes can therefore be reduced to

1015    minimize within haploblock recombination and better define QTL regions. However, when haploblocks

1016    are very small, many haploblocks will consist of only one SNP or a few SNPs, increasing data sizes (and

1017    thereby computation time in downstream analyses) and reducing the number of haplotypes per

1018    haploblock, which can reduce the suitability of the data for visual examination. Unlike the 8K apple SNP

1019    array, the 20K apple SNP array was designed to have clusters of multiple SNPs spread at approximately 1

1020    cM intervals. A similar approach was used to create 9K add-ons for the 9K peach array and 6K cherry

1021    array [59]. This strategy supports the generation of haploblocks consisting of SNPs aggregated within 1

1022    cM intervals while still having multiple SNPs in a single haploblock and thus multiple informative

1023    haplotypes.

1024

1025    Different germplasm will also lead to different haploblock borders. Currently, haploblocks are based on

1026    historic recombination events representing the U.S. breeding programs included in this study. Other

1027    breeding programs or genetic studies might have other sets of founders and thus different

1028    recombinations of relevance. Furthermore, the addition of new advanced selections and parents will

1029    introduce new recombinations in their germplasm. Finally, as the understanding of the apple, peach,

1030    and cherry germplasm increases, previously unknown progenitors, founders, and pedigree connections

1031    will be discovered, also increasing the number of observed recombinations.

1032

1033    Given that haploblocking is performed at a relative late stage in the workflow, haploblock borders can

1034    be altered without the need to redo all previously conducted pedigree and SNP marker curation. In fact,

1035    existing haplotype data can be converted back to phased, fully curated SNP data which, in turn, can be

1036    used to determine haplotypes for any set of haploblocks. As the SNPs are already phased and missing

1037    SNP data was imputed based on the haplotypes, haplotype determination for new haploblock borders

1038    should not create new genotyping errors in the data set. Once numbers of new recombinations are high

1039    enough to justify updating of haploblock data, part of the haploblocks and their haplotypes should be

1040    altered. PediHaplotyper supports the use of previous haplotype definitions for haploblocks that did not

1041    change in composition. Adjusted haploblocks could be marked through their names, thus providing tools

1042    to monitor new as well as previous, possibly well-known, marker alleles.

1043

1044    *The SNP classification system and integration of genotypic data for new germplasm into existing data*

1045    *sets*

1046    The established SNP classification system enables the quick creation of a subset of SNPs that require

1047    minimal or no data curation and provides a guideline on possible issues with other SNPs and how to

1048    solve them. The system should help with the quick integration of new genotypic data into existing data

1049    sets. Genotype calls for SNPs of type 1 and type 2 can be quickly integrated with high confidence in their

1050    genotype calls. Where desired, SNPs of type 3, 4, and 5 can also be integrated, but additional curation

46

1051    would be required. Depending on germplasm tested, these SNPs might have incorrect genotype scoring

1052    but their SNP type is an indication of why the genotype scoring is wrong and how to fix it. In other

1053    germplasm, additional SNPs in the probe or null alleles might not be present, causing SNPs that are now

1054    classified as type 4 or type 5 to give reliable results as if they were type 1 or type 2. Similarly, if

1055    germplasm is used that is unrelated to that used here, type 1 and type 2 SNPs might show additional

1056    clusters or null alleles and will require further curation. Finally, type 7 SNPs, which could not be mapped

1057    in this germplasm, might be mapped and valuable for other germplasm.

1058

1059    The available reference data (www.rosaceae.org), combined with the SNP classification system, will

1060    facilitate correct curation of additional genotypic data, even if the new germplasm is not directly

1061    descended. The SNP genotype calls provided here are a reference for the genotype of each observed

1062    genotype cluster in GenomeStudio®. In addition, SNP cluster coordinates of the latest GenomeStudio®

1063    file can be imported into new projects, thus helping GenomeStudio® to correctly identify clusters.

1064    Finally, the use of reference iScan data is especially useful for markers that have only two of the three

1065    clusters in a new project but all three clusters were defined in the current reference dataset. By adding

1066    reference iScan data into the new project, all three clusters will be available, ensuring correct

1067    automated genotype calling. Therefore, we recommend including available reference data when

1068    obtaining genotype calls for new germplasm.

1069

1070    *Data curation in apple*

1071    The need for SNP data curation in apple was increased by the whole genome duplication in the

1072    evolutionary history of apple, the relatively poorer quality of the first genome draft used for

1073    development of the 8K SNP array, and unidentified polymorphisms in the probe regions during SNP

1074    array design. The genome duplication resulted in presence of multiple highly similar sequences on

47

1075    different chromosomes. Indeed, a BLAST analysis against *Malus* genome v1.0 of the first 24 nucleotides

1076    of the 3' region of arrayed SNP probes, which is most important for probe binding, showed that

1077    approximately 50% of the sequences returned multiple hits with almost all of these hits being located on

1078    multiple LGs [33]. This proportion is expected to be lower for the latest genome version [44] as most

1079    errors in assembly were removed but the proportion is expected to remain high due to chromosome

1080    and gene duplication observed in apple. Where two genomic regions are targeted by the same probe,

1081    complex cluster plots will occur if more than one of the targeted loci segregate within a single family.

1082    Such markers must be excluded from a curated data set. Multi-target markers might still be robust if

1083    they segregate at only one locus. In this case, only the cluster plot space is reduced (mostly halved),

1084    causing clusters to be located more closely to each other. In turn, this might occasionally cause

1085    separation issues. Also, some markers are lost because GenomeStudio® cannot assign genotype calls for

1086    markers where one of the homozygous clusters is located at theta=0.5, the center of the x-axis, and thus

1087    these markers are considered by the software to have failed. . A special case for two-locus markers

1088    occurs where each locus segregates in specific families but both loci never segregate together in the

1089    same family. In this case, genotype scoring might be performed accurately, and the SNP still needs to be

1090    present twice in the map although under different names. Two- and three-locus SNPs have been

1091    successfully mapped in the multi-family based genetic linkage map created by Di Pierro and co-workers

1092    (2016) [45]. However, in subsequent QTL mapping studies on pedigreed germplasm, such markers were

1093    excluded, as in the current study.

1094

1095    Several intermediate progenitors in the apple data set lacked any genotypic data and therefore the

1096    recorded link between some important breeding parents and their ancestors had to be set to unknown

1097    during haploblock and haplotype determination. For some progenitors, 20K data from the European

1098    FruitBreedomics project was available that reestablished the connection between genotyped individuals

48

1099    and their ancestors, but many other progenitors likely no longer exist. Individuals that were

1100    disconnected from the pedigree with little representation could not therefore have their haplotypes

1101    accurately determined using PediHaplotyper. It was, however, possible to manually determine their

1102    haplotypes based on their SNP data and haplotypes present in disconnected relatives.

1103

1104    *Data curation in peach*

1105    In peach, the most challenging step in the workflow was the curation of pedigree information over nine

1106    generations. Although much pedigree information is available in the literature [60], we identified

1107    incorrect parentage in the PC error analysis in cultivars and breeding selections, which we attributed to

1108    selfing or outcrossing. Incorrect pedigree records were previously reported in the UC Davis processing

1109    peach breeding program in approximately 20% of individuals, both parental and breeding selections

1110    [16]. In this work, we identified incorrect parentage in approximately 11% of the pedigree records from

1111    the three fresh market peach breeding programs, most of which were observed in breeding selections.

1112    High level of inbreeding and coancestry in the U.S cultivated peach germplasm [61] creates overlap in

1113    the ancestral generations of most U.S. peach breeding programs. Therefore, corrections in the ancestral

1114    pedigree records reported by Fresnedo-Ramírez and co-workers (2015) [16] reduced the number of

1115    errors detected here. Furthermore, intermediate parents were unavailable for genotyping, so pedigree

1116    connections were preserved by retaining pedigree information even though many intermediate

1117    progenitors were not genotyped. Finally, the presence of missing data within a haplotype resulted in

1118    Mendelian-inconsistent errors in the haploblock and haplotype generation steps, which made the

1119    haploblock data curation time-consuming.

1120

1121

1122

1123    *Data curation in sweet cherry*

1124    For the sweet cherry germplasm, the most challenging issue was the small sample size of some families

1125    (as few as four individuals), which were too small for FlexQTL™ to accurately determine linkage phase.

1126    For parents with just one genotyped offspring, phasing of the parent homologs was considered putative

1127    as recombination inherited by offspring could not be determined. For those parents with just two

1128    genotyped offspring, recombinations were arbitrarily assigned between the two offspring, as the true

1129    recombinant offspring could not be determined. In addition, scarce information on pedigrees in

1130    ancestral generations beyond about five limited further imputations in data curation, unlike for apple

1131    and peach. Various founders showed extensive regions of common haplotypes, indicating a high degree

1132    of relatedness among such founders. Some recently published haplotyping results exemplify this for the

1133    founders 'Black Republican' and 'Napoleon' [21]. Unraveling the unknown relationships among founders

1134    could thus provide useful information for future data curation in sweet cherry.

1135

1136    *Expectations for other crops*

1137    The proposed workflow could be applied to other diploid crops with similar breeding systems where

1138    clonally propagated relatives of current breeding material still exist. However, there are additional

1139    aspects that would need to be considered in certain circumstances that were not encountered in the

1140    present study. First, this workflow makes the assumption that there are no differences in the true SNP

1141    map order among individuals of a species. In interspecific crosses where there can be differences in

1142    chromosome arrangements between parental species, the different SNP order or indel variation among

1143    individuals could result in additional perceived double recombinations or other difficulties in following

1144    this workflow. Additionally, this workflow assumes that there is sufficient marker information to

1145    correctly identify pedigree relationships and assumes sufficient segregation information for validating

1146    marker order and identifying Mendelian-consistent errors. When using highly homozygous, inbred

50

1147   individuals, there might be too few segregating markers available to correctly identify marker order or

1148   find Mendelian-consistent errors through double recombinations. Also, for small germplasm sets, too

1149   few recombinations might be available to detect incorrect marker order. Finally, the prevalence of

1150   missing genotypic values should be sufficiently low across individuals. Unlike the SNP arrays used in this

1151   study, some genotyping methods such as Genotyping-by-Sequencing do not consistently target specific

1152   loci. This non-specificity can increase the flexibility of their use, but also raises new issues for which the

1153   current workflow would have to be adapted, including the potential decrease in accuracies of

1154   genotyping and haploblock determination due to unbalanced representation of genotyped loci, high

1155   levels of missing data, and sequencing errors.

1156

1157   *High-quality archived SNP and haplotype data sets*

1158   The presented genome-wide genotypic data sets for apple, peach, and sweet cherry are of very high

1159   quality, are composed of genetically complex germplasm, and contain no errors that could be

1160   determined based on pedigree information. This high quality provides confidence in the results of

1161   downstream analyses. Such confidence is important as many of these results are expected to lead to

1162   fundamental discoveries and practical breeding application. The iScan data, phased SNP, and haplotype

1163   datasets of individuals in the apple, peach, and sweet cherry crop reference sets are available through

1164   the Genome Database for Rosaceae (www.rosaceae.org).

1165

1166   Marker and pedigree data from germplasm subsets of the current U.S. RosBREED project, the EU-

1167   FruitBreedomics project, and other research projects have previously been curated by a precursor to the

1168   current workflow and used for the creation of a multi-family based genetic linkage maps [20,45] and in

1169   multifamily based QTL studies in apple [62–65], peach [22,66], and sweet cherry [14]. Also, elements of

1170   the workflow were used for allo-octoploid strawberry to curate Axiom-based SNP markers [31] and

51

1171    pedigree data that were subsequently used in multi-family based QTL analyses [67–69]. While providing

1172    high-quality data for each analysis separately, these earlier steps in data curation have helped guide and

1173    streamline the data curation workflow presented here. The current workflow and resulting data sets

1174    ensure that the same curation steps have been used across the data sets of multiple crops and that the

1175    data sets are of the same high quality.

1176

1177    **Conclusion**

1178    A curation workflow for genotypic data of pedigreed germplasm was generated by determining the

1179    optimal order of resolving issues and by providing a step-by-step guideline. Using simple genetic

1180    principles, errors can be found and curated in a directed and efficient way, reducing the time needed to

1181    obtain a high-quality genotypic data set. The workflow was used to obtain a SNP data set for large

1182    germplasm sets for each of apple, peach, and sweet cherry representing U.S. breeding programs based

1183    on the apple 8K SNP array, peach 9K SNP array, and cherry 6K SNP array, respectively, whose SNP data is

1184    available through this paper (www.rosaceae.org), as well as used on apple and peach germplasm sets

1185    representing European breeding programs based on the apple 20K and peach 9K arrays, whose SNP data

1186    are still private. These high-quality data sets contain the largest sets of SNPs obtained through their

1187    respective SNP arrays and will provide the foundation for confident subsequent analyses in genetic

1188    research.

1189

1190    **Acknowledgements**

1195    USDA NIFA Hatch projects 0211277and 1014919, and the FruitBreedomics project No 265582:

1196    Integrated approach for increasing breeding efficiency in fruit tree crops (www.FruitBreedomics.com)

1197    that was co-funded by the EU seventh Framework Programme.

1198

1199    **Supplementary information**

1200    **Table S1:** Apple germplasm genotyped and used for data curation workflow. Individuals are split over

1201    the publicly available RosBREED Crop Reference Set, three privately held RosBREED Breeding Pedigree

1202    Sets, and genotypic data received from either KULeuven (Belgium) or the FruitBreedomics project.

1203    Except for the Breeding Pedigree sets, curated pedigree information is given for each individual. For

1204    each individual, the type of material (selected vs. unselected), the location of sampling, quality of the

1205    results, and inferred ploidy of the sample are given. For unselected seedlings, the family to which they

1206    belong is also given. For the Breeding Pedigree Sets, this information is summarized per full-sib family. If

1207    tissue was collected at the USDA germplasm repository in Geneva, a GRIN accession number is also

1208    provided. Parents highlighted in yellow did not have genotypic data and their pedigree-relationships

1209    could not be tested.

1210

1211    **Table S2:** Peach germplasm genotyped and used for curation workflow. Individuals are split over the

1212    publicly available RosBREED Crop Reference Set and three privately held RosBREED Breeding Pedigree

1213    Sets. Except for the Breeding Pedigree Sets, curated pedigree information is given for each individual.

1214    For each individual, the type of material (selected vs. unselected), the location of sampling, and quality

1215    of the results of the sample are given. For unselected seedlings, the family to which they belong is also

1216    given. For the Breeding Pedigree Sets, this information is summarized per full-sib family.

1217

53

1218    **Table S3:** Sweet cherry germplasm genotyped and used for curation workflow. All individuals are part of

1219    the publicly available RosBREED Crop Reference Set. For each individual, curated pedigree information,

1220    the type of material (selected vs. unselected), the location of sampling, and quality of the results of the

1221    sample are given. For unselected seedlings, the family to which they belong is also given.

1222

1223    **Table S4:** Parameter settings used for (A) filtering SNPs used in analyses of B-allele frequency, (B)

1224    running ASSiST, (C) running FlexQTL™ for detecting Mendelian-inconsistent errors and Mendelian-

1225    consistent errors, and (D) running FlexQTL™ for phasing, haploblock determination, and creating

1226    PediHaplotyper input files.

1227

1228    **Table S5:** Final genetic map used for apple during data curation. For each marker, genetic position,

1229    associated haploblock, and physical position based on the apple GDDH 13 v1.1 genome are given.

1230

1231    **Table S6:** Final genetic map used for peach during data curation. For each marker, genetic position,

1232    associated haploblock, and physical position based on the peach v2 genome are given.

1233

1234    **Table S7:** Final genetic map used for sweet cherry during data curation. For each marker, genetic

1235    position, associated haploblock, and physical position based on the peach v2 genome are given.

1236

1237    **Table S8:** SNP classification for apple. Each SNP is classified as follows: Type '1' for SNPs with good

1238    clustering and less than 5% call errors, '2' for SNPs with shifted clusters causing one of the clusters to be

1239    called incorrectly, '3' for SNPs with additional clusters (excluding null-alleles) that cause the incorrect

1240    identification of at least one cluster, '4' for SNPs with null-alleles that cannot be correctly called

1241    automatically, '5' for SNPs that could not be mapped accurately but had correct clustering, '6' for

1242    monomorphic SNPs, and '7' for failed SNPs.

1243

1244    **Table S9:** SNP classification for peach. Each SNP is classified as follows: Type '1' for SNPs with good

1245    clustering and less than 5% call errors, '2' for SNPs with shifted clusters causing one of the clusters to be

1246    called incorrectly, '3' for SNPs with additional clusters (excluding null-alleles) that cause the incorrect

1247    identification of at least one cluster, '4' for SNPs with null-alleles that cannot be correctly called

1248    automatically, '5' for SNPs that could not be mapped accurately but had correct clustering, '6' for

1249    monomorphic SNPs, and '7' for failed SNPs.

1250

1251    **Table S10:** SNP classification for sweet cherry. Each SNP is classified as follows: Type '1' for SNPs with

1252    good clustering and less than 5% call errors, '2' for SNPs with shifted clusters causing one of the clusters

1253    to be called incorrectly, '3' for SNPs with additional clusters (excluding null-alleles) that cause the

1254    incorrect identification of at least one cluster, '4' for SNPs with null-alleles that cannot be correctly

1255    called automatically, '5' for SNPs that could not be mapped accurately but had correct clustering, '6' for

1256    monomorphic SNPs, and '7' for failed SNPs.

1257

1258    **Figure S1:** SNP B-allele frequences plotted against physical position in the genome for (A) triploid

1259    individuals excluding 'Jonagold', and (B) individuals with a tetraploid pattern

1260

1261    **Figure S2:** Call rates observed for individuals classified as having good, intermediate, or bad quality of

1262    genotypic data as defined by their B-allele frequency plot outcome. Higher call rates are observed for

1263    individuals with better quality of genotypic data.

1264

1265    **Document S1:** R-script used to create B-allele frequency plots for all genotyped individuals.

1266

1267    **Document S2:** R-scripts used to confirm and deduce P(P)C relationships.

1268

1269    **Document S3:** Hands-on guideline on how to perform data curation using the steps described in this

1270    study

1271

1272    **References**

1273    1.    Pompanon F, Bonin A, Bellemain E, Taberlet P. Genotyping errors: causes, consequences and
1274          solutions. Nat Rev Genet. 2005;6: 847–859. doi:10.1038/nrg1707

1275    2.    Buetow KH. Influence of aberrant observations on high-resolution linkage analysis outcomes. Am J
1276          Hum Genet. 1991;49: 985–994.

1277    3.    Goldstein DR, Zhao H, Speed TP. The effects of genotyping errors and interference on estimation of
1278          genetic distance. Hum Hered. 1997;47: 86–100. doi:10.1159/000154396

1279    4.    Hackett CA, Broadfoot LB. Effects of genotyping errors, missing values and segregation distortion
1280          in molecular marker data on the construction of linkage maps. Heredity. 2003;90: 33–38.
1281          doi:10.1038/sj.hdy.6800173

1282    5.    Abecasis GR, Cherny SS, Cardon LR. The impact of genotyping error on family-based analysis of
1283          quantitative traits. European Journal of Human Genetics. 2001;9: 130–134.
1284          doi:10.1038/sj.ejhg.5200594

1285    6.    Gordon D, Finch SJ, Nothnagel M, Ott J. Power and sample size calculations for case-control
1286          genetic association tests when errors are present: application to single nucleotide polymorphisms.
1287          Hum Hered. 2002;54: 22–33. doi:10.1159/000066696

1288    7.    Browning BL, Yu Z. Simultaneous genotype calling and haplotype phasing improves genotype
1289          accuracy and reduces false-positive associations for genome-wide association studies. Am J Hum
1290          Genet. 2009;85: 847–861. doi:10.1016/j.ajhg.2009.11.004

1291    8.    Terwilliger J, Weeks D, Ott J. Laboratory errors in the reading of marker alleles cause massive
1292          reductions in LOD score and lead to gross overestimates of the recombination fraction. Am J Hum
1293          Genet. 1990;47: A201.

1294    9.    Kirk KM, Cardon LR. The impact of genotyping error on haplotype reconstruction and frequency
1295          estimation. Eur J Hum Genet. 2002;10: 616–622. doi:10.1038/sj.ejhg.5200855

1296  10.  Cheung CYK, Thompson EA, Wijsman EM. Detection of Mendelian consistent genotyping errors in
1297       pedigrees. Genet Epidemiol. 2014;38: 291–299. doi:10.1002/gepi.21806

1298  11.  Vouillamoz JF, Grando MS. Genealogy of wine grape cultivars: "Pinot" is related to "Syrah."
1299       Heredity (Edinb). 2006;97: 102–110. doi:10.1038/sj.hdy.6800842

1300  12.  Evans KM, Patocchi A, Rezzonico F, Mathis F, Durel CE, Fernández-Fernández F, et al. Genotyping
1301       of pedigreed apple breeding material with a genome-covering set of SSRs: trueness-to-type of
1302       cultivars and their parentages. Mol Breeding. 2011;28: 535–547. doi:10.1007/s11032-010-9502-5

1303  13.  Lacombe T, Boursiquot J-M, Laucou V, Di Vecchi-Staraz M, Péros J-P, This P. Large-scale parentage
1304       analysis in an extended set of grapevine cultivars (Vitis vinifera L.). Theor Appl Genet. 2013;126:
1305       401–414. doi:10.1007/s00122-012-1988-2

1306  14.  Rosyara UR, Sebolt AM, Peace C, Iezzoni AF. Identification of the paternal parent of 'Bing' sweet
1307       cherry and confirmation of descendants using single nucleotide polymorphism markers. J Amer Soc
1308       Hort Sci. 2014;139: 148–156.

1309  15.  Pikunova A, Madduri M, Sedov E, Noordijk Y, Peil A, Troggio M, et al. 'Schmidt's Antonovka' is
1310       identical to 'Common Antonovka', an apple cultivar widely used in Russia in breeding for biotic and
1311       abiotic stresses. Tree Genetics & Genomes. 2014;10: 261–271. doi:10.1007/s11295-013-0679-8

1312  16.  Fresnedo-Ramírez J, Crisosto CH, Gradziel TM, Famula TR. Pedigree correction and estimation of
1313       breeding values for peach genetic improvement. Acta Horticulturae. 2015;1084: 249–256.
1314       doi:10.17660/ActaHortic.2015.1084.35

1315  17.  Fresnedo-Ramírez J, Frett TJ, Sandefur PJ, Salgado-Rojas A, Clark JR, Gasic K, et al. QTL mapping
1316       and breeding value estimation through pedigree-based analysis of fruit size and weight in four
1317       diverse peach breeding programs. Tree Genetics & Genomes. 2016;12: 25. doi:10.1007/s11295-
1318       016-0985-z

1319  18.  Lassois L, Denancé C, Ravon E, Guyader A, Guisnel R, Hibrand-Saint-Oyant L, et al. Genetic
1320       diversity, population structure, parentage analysis, and construction of core collections in the
1321       French apple germplasm based on SSR markers. Plant Mol Biol Rep. 2016;34: 827–844.
1322       doi:10.1007/s11105-015-0966-7

1323  19.  Larsen B, Toldam-Andersen TB, Pedersen C, Ørgaard M. Unravelling genetic diversity and cultivar
1324       parentage in the Danish apple gene bank collection. Tree Genetics & Genomes. 2017;13: 14.
1325       doi:10.1007/s11295-016-1087-7

1326  20.  Howard NP, van de Weg E, Bedford DS, Peace CP, Vanderzande S, Clark MD, et al. Elucidation of
1327       the 'Honeycrisp' pedigree through haplotype analysis with a multi-family integrated SNP linkage
1328       map and a large apple (*Malus×domestica*) pedigree-connected SNP data set. Horticulture
1329       Research. 2017;4: 17003. doi:10.1038/hortres.2017.3

1330  21.  Cai L, Voorrips RE, van de Weg E, Peace C, Iezzoni A. Genetic structure of a QTL hotspot on
1331       chromosome 2 in sweet cherry indicates positive selection for favorable haplotypes. Mol Breeding.
1332       2017;37: 85. doi:10.1007/s11032-017-0689-6

1333   22.   Hernández Mora JR, Micheletti D, Bink MAM, Van de Weg WE, Bassi D, Nazzicari N, et al.
1334         Discovering peach QTLs with multiple progeny analysis. Acta Horticulturae. 2017;1172: 405–410.
1335         doi:10.17660/ActaHortic.2017.1172.77

1336   23.   van de Weg E, Di Guardo M, Jänsch M, Socquet-Juglard D, Costa F, Baumgartner I, et al. Epistatic
1337         fire blight resistance QTL alleles in the apple cultivar 'Enterprise' and selection X-6398 discovered
1338         and characterized through pedigree-informed analysis. Mol Breeding. 2018;38: 5.
1339         doi:10.1007/s11032-017-0755-0

1340   24.   Xu X, Bai G. Whole-genome resequencing: changing the paradigms of SNP detection, molecular
1341         mapping and gene discovery. Mol Breeding. 2015;35: 33. doi:10.1007/s11032-015-0240-6

1342   25.   Rasheed A, Hao Y, Xia X, Khan A, Xu Y, Varshney RK, et al. Crop breeding chips and genotyping
1343         platforms: progress, challenges, and perspectives. Mol Plant. 2017;10: 1047–1064.
1344         doi:10.1016/j.molp.2017.06.008

1345   26.   Chagné D, Crowhurst RN, Troggio M, Davey MW, Gilmore B, Lawley C, et al. Genome-wide SNP
1346         detection, validation, and development of an 8K SNP array for apple. PLOS ONE. 2012;7: e31745.
1347         doi:10.1371/journal.pone.0031745

1348   27.   Verde I, Bassil N, Scalabrin S, Gilmore B, Lawley CT, Gasic K, et al. Development and evaluation of a
1349         9K SNP array for peach by internationally coordinated SNP detection and validation in breeding
1350         germplasm. PLOS ONE. 2012;7: e35668. doi:10.1371/journal.pone.0035668

1351   28.   Peace C, Bassil N, Main D, Ficklin S, Rosyara UR, Stegmeir T, et al. Development and evaluation of a
1352         genome-wide 6K SNP array for diploid sweet cherry and tetraploid sour cherry. PLOS ONE. 2012;7:
1353         e48305. doi:10.1371/journal.pone.0048305

1354   29.   Le Paslier M-C, Choisne N, Bacilieri R, Bounon R, Boursiquot J-M, Bras M, et al. The GrapeReSeq
1355         18K Vitis genotyping chip. 9th International symposium grapevine physiology and biotechnology.
1356         La Serena, Chili; 2013. p. 123.

1357   30.   Bianco L, Cestaro A, Sargent DJ, Banchi E, Derdak S, Guardo MD, et al. Development and validation
1358         of a 20K single nucleotide polymorphism (SNP) whole genome genotyping array for apple (*Malus ×
1359         domestica* Borkh). PLOS ONE. 2014;9: e110377. doi:10.1371/journal.pone.0110377

1360   31.   Bassil NV, Davis TM, Zhang H, Ficklin S, Mittmann M, Webster T, et al. Development and
1361         preliminary evaluation of a 90K Axiom® SNP array for the allo-octoploid cultivated strawberry
1362         *Fragaria× ananassa*. BMC Genomics. 2015;16: 155. doi:10.1186/s12864-015-1310-1

1363   32.   Bianco L, Cestaro A, Linsmith G, Muranty H, Denancé C, Théron A, et al. Development and
1364         validation of the Axiom® Apple 480K SNP genotyping array. Plant J. 2016;86: 62–74.
1365         doi:10.1111/tpj.13145

1366   33.   Troggio M, Šurbanovski N, Bianco L, Moretto M, Giongo L, Banchi E, et al. Evaluation of SNP data
1367         from the *Malus* Infinium array identifies challenges for genetic analysis of complex genomes of
1368         polyploid origin. PLOS ONE. 2013;8: e67407. doi:10.1371/journal.pone.0067407

1369    34.    Di Guardo M, Micheletti D, Bianco L, Koehorst-van Putten HJJ, Longhi S, Costa F, et al. ASSIsT: an
1370            automatic SNP scoring tool for in- and outbreeding species - Reference Manual. 2015.

1371    35.    Di Guardo M, Micheletti D, Bianco L, Koehorst-van Putten HJJ, Longhi S, Costa F, et al. ASSIsT: an
1372            automatic SNP scoring tool for in- and outbreeding species. Bioinformatics. 2015;31: 3873–3874.
1373            doi:10.1093/bioinformatics/btv446

1374    36.    Voorrips RE, Bink MCAM, Kruisselbrink JW, Koehorst-van Putten HJJ, van de Weg WE.
1375            PediHaplotyper: software for consistent assignment of marker haplotypes in pedigrees. Mol Breed.
1376            2016;36. doi:10.1007/s11032-016-0539-y

1377    37.    Iezzoni A, Weebadde C, Luby J, Chengyan Yue, van de Weg E, Fazio G, et al. Rosbreed: enabling
1378            marker-assisted breeding in Rosaceae. Acta Horticulturae. 2010;859: 389–394.
1379            doi:10.17660/ActaHortic.2010.859.47

1380    38.    Iezzoni A, Weebadde C, Peace C, Main D, Bassil NV, Coe M, et al. Where are we now as we merge
1381            genomics into plant breeding and what are our limitations? Experiences from RosBREED. Acta
1382            Horticulturae. 2016;1117: 1–5. doi:10.17660/ActaHortic.2016.1117.1

1383    39.    Iezzoni A, Peace C, Main D, Bassil N, Coe M, Finn C, et al. RosBREED2: progress and future plans to
1384            enable DNA-informed breeding in the *Rosaceae*. Acta Horticulturae. 2017;1172: 115–118.
1385            doi:10.17660/ActaHortic.2017.1172.20

1386    40.    Laurens F, Durel C-E, Patocchi A, Peil A, Salvi S, Tartarini S, et al. Review on apple genetics and
1387            breeding programs and presentation of a new initiative of a news European initiative to increase
1388            fruit breeding efficiency. Journal of Fruit Science. 2010;27: 102–107.

1389    41.    Laurens F, Aranzana MJ, Arús P, Bassi D, Bonany J, Corelli L, et al. Review of fruit genetics and
1390            breeding programmes and a new European initiative to increase fruit breeding efficiency. Acta
1391            Horticulturae. 2012;929: 95–102. doi:10.17660/ActaHortic.2012.929.12

1392    42.    Laurens F, Aranzana MJ, Arus P, Bassi D, Bink M, Bonany J, et al. An integrated approach for
1393            increasing breeding efficiency in apple and peach in Europe. Hortic Res. 2018;5: 11.
1394            doi:10.1038/s41438-018-0016-3

1395    43.    Peace CP, Luby JJ, van de Weg WE, Bink MCAM, Iezzoni AF. A strategy for developing
1396            representative germplasm sets for systematic QTL validation, demonstrated for apple, peach, and
1397            sweet cherry. Tree Genetics & Genomes. 2014;10: 1679–1694. doi:10.1007/s11295-014-0788-z

1398    44.    Daccord N, Celton J-M, Linsmith G, Becker C, Choisne N, Schijlen E, et al. High-quality de novo
1399            assembly of the apple genome and methylome dynamics of early fruit development. Nat Genet.
1400            2017;49: 1099–1106. doi:10.1038/ng.3886

1401    45.    Di Pierro EA, Gianfranceschi L, Di Guardo M, Koehorst-van Putten HJ, Kruisselbrink JW, Longhi S, et
1402            al. A high-density, multi-parental SNP genetic map on apple validates a new mapping approach for
1403            outcrossing species. Horticulture Research. 2016;3: 16057. doi:10.1038/hortres.2016.57

1404 46. Verde I, Jenkins J, Dondini L, Micali S, Pagliarani G, Vendramin E, et al. The Peach v2.0 release:
1405      high-resolution linkage mapping and deep resequencing improve chromosome-scale assembly and
1406      contiguity. BMC Genomics. 2017;18: 225. doi:10.1186/s12864-017-3606-9

1407 47. Klagges C, Campoy JA, Quero-García J, Guzmán A, Mansur L, Gratacós E, et al. Construction and
1408      comparative analyses of highly dense linkage maps of two sweet cherry intra-specific progenies of
1409      commercial cultivars. PLOS ONE. 2013;8: e54743. doi:10.1371/journal.pone.0054743

1410 48. Chagné D, Kirk C, Whitworth C, Erasmuson S, Bicknell R, Sargent DJ, et al. Polyploid and aneuploid
1411      detection in apple using a single nucleotide polymorphism array. Tree Genetics & Genomes.
1412      2015;11: 94. doi:10.1007/s11295-015-0920-8

1413 49. R Core Team. R: A Language and Environment for Statistical Computing [Internet]. Vienna, Austria:
1414      R Foundation for Statistical Computing; 2018. Available: http://www.R-project.org/

1415 50. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: a tool set for
1416      whole-genome association and population-based linkage analyses. Am J Hum Genet. 2007;81:
1417      559–575. doi:10.1086/519795

1418 51. Kalinowski ST, Taper ML, Marshall TC. Revising how the computer program CERVUS
1419      accommodates genotyping error increases success in paternity assignment. Mol Ecol. 2007;16:
1420      1099–1106. doi:10.1111/j.1365-294X.2007.03089.x

1421 52. Young ND, Tanksley SD. Restriction fragment length polymorphism maps and the concept of
1422      graphical genotypes. Theoret Appl Genetics. 1989;77: 95–101. doi:10.1007/BF00292322

1423 53. Namjou B, Sestak AL, Armstrong DL, Zidovetzki R, Kelly JA, Jacob N, et al. High-density genotyping
1424      of *STAT4* reveals multiple haplotypic associations with systemic lupus erythematosus in different
1425      racial groups. Arthritis Rheum. 2009;60: 1085–1095. doi:10.1002/art.24387

1426 54. Jacob CO, Zhu J, Armstrong DL, Yan M, Han J, Zhou XJ, et al. Identification of *IRAK1* as a risk gene
1427      with critical role in the pathogenesis of systemic lupus erythematosus. Proc Natl Acad Sci USA.
1428      2009;106: 6256–6261. doi:10.1073/pnas.0901181106

1429 55. Kumar S, Chagné D, Bink MCAM, Volz RK, Whitworth C, Carlisle C. Genomic selection for fruit
1430      quality traits in apple (Malus×domestica Borkh.). PLOS ONE. 2012;7: e36674.
1431      doi:10.1371/journal.pone.0036674

1432 56. Guo Y, He J, Zhao S, Wu H, Zhong X, Sheng Q, et al. Illumina human exome genotyping array
1433      clustering and quality control. Nat Protoc. 2014;9: 2643–2662. doi:10.1038/nprot.2014.174

1434 57. Vanderzande S, Micheletti D, Troggio M, Davey MW, Keulemans J. Genetic diversity, population
1435      structure, and linkage disequilibrium of elite and local apple accessions from Belgium using the
1436      IRSC array. Tree Genetics & Genomes. 2017;13: 125. doi:10.1007/s11295-017-1206-0

1437 58. da Silva Linge C, Antanaviciute L, Abdelghafar A, Arús P, Bassi D, Rossini L, et al. High-density multi-
1438      population consensus genetic linkage map for peach. PLOS ONE. 2018;13: e0207724.
1439      doi:10.1371/journal.pone.0207724

1440    59.    Vanderzande S, Zheng P, Cai L, Iezzoni A, Main D, Peace C. Development and initial assessment of
1441        the 9K SNP addition to the sweet and sour cherry genome-wide SNP array. San Diego, CA, USA;
1442        2019.

1443    60.    Okie WR. Handbook of peach and nectarine varieties: performance in the southeastern United
1444        States and index of names. U.S. Dept. of Agriculture, Agricultural Research Service; 1998.

1445    61.    Scorza R, Sherman W. Peaches. In: Janick J, Moore J, editors. Fruit Breeding. New York: Wiley;
1446        1996. pp. 325–440.

1447    62.    Allard A, Bink MCAM, Martinez S, Kelner J-J, Legave J-M, di Guardo M, et al. Detecting QTLs and
1448        putative candidate genes involved in budbreak and flowering time in an apple multiparental
1449        population. J Exp Bot. 2016;67: 2875–2888. doi:10.1093/jxb/erw130

1450    63.    Di Guardo M, Bink MCAM, Guerra W, Letschka T, Lozano L, Busatto N, et al. Deciphering the
1451        genetic control of fruit texture in apple by multiple family-based analysis and genome-wide
1452        association. J Exp Bot. 2017;68: 1451–1466. doi:10.1093/jxb/erx017

1453    64.    Durand J-B, Allard A, Guitton B, van de Weg E, Bink MCAM, Costes E. Predicting flowering behavior
1454        and exploring its genetic determinism in an apple multi-family population based on statistical
1455        indices and simplified phenotyping. Front Plant Sci. 2017;8: 858. doi:10.3389/fpls.2017.00858

1456    65.    Verma S, Evans K, Guan Y, Luby JJ, Rosyara UR, Howard NP, et al. Two large-effect QTLs, Ma and
1457        Ma3, determine genetic potential for acidity in apple fruit: Breeding insights from a multi-family
1458        study. Tree Genetics & Genomes. 2019; *in press*.

1459    66.    Hernández Mora JR, Micheletti D, Bink M, Van de Weg E, Cantín C, Nazzicari N, et al. Integrated
1460        QTL detection for key breeding traits in multiple peach progenies. BMC Genomics. 2017;18: 404.
1461        doi:10.1186/s12864-017-3783-6

1462    67.    Roach JA, Verma S, Peres NA, Jamieson AR, van de Weg WE, Bink MCAM, et al. FaRXf1: a locus
1463        conferring resistance to angular leaf spot caused by *Xanthomonas fragariae* in octoploid
1464        strawberry. Theor Appl Genet. 2016;129: 1191–1201. doi:10.1007/s00122-016-2695-1

1465    68.    Mangandi J, Verma S, Osorio L, Peres NA, Weg E van de, Whitaker VM. Pedigree-Based Analysis in
1466        a multiparental population of octoploid strawberry reveals QTL alleles conferring resistance to
1467        *Phytophthora cactorum*. G3: Genes, Genomes, Genetics. 2017;7: 1707–1719.
1468        doi:10.1534/g3.117.042119

1469    69.    Anciro A, Mangandi J, Verma S, Peres N, Whitaker VM, Lee S. FaRCg1: a quantitative trait locus
1470        conferring resistance to Colletotrichum crown rot caused by *Colletotrichum gloeosporioides* in
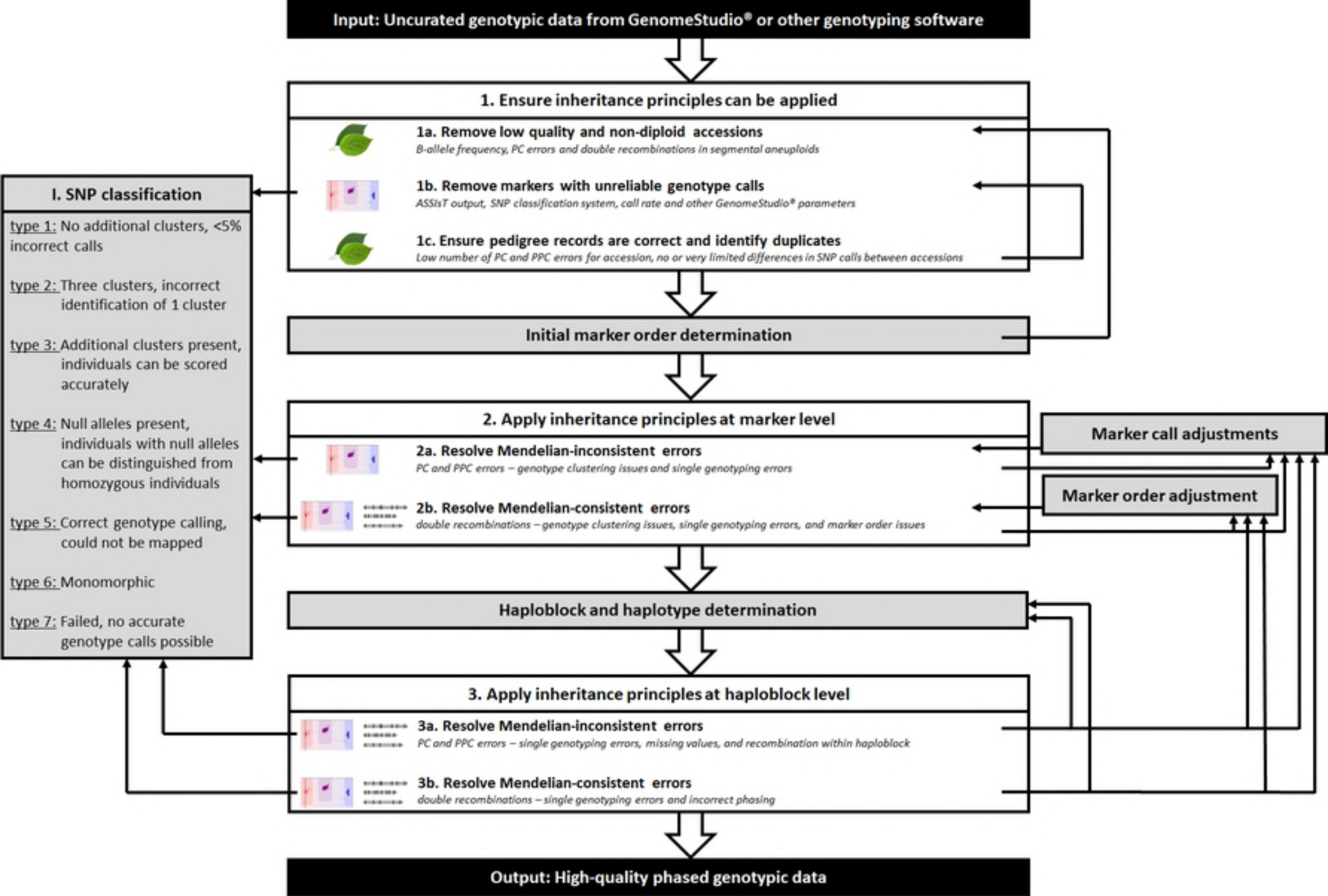1471        octoploid strawberry. Theor Appl Genet. 2018;131: 2167–2177. doi:10.1007/s00122-018-3145-z
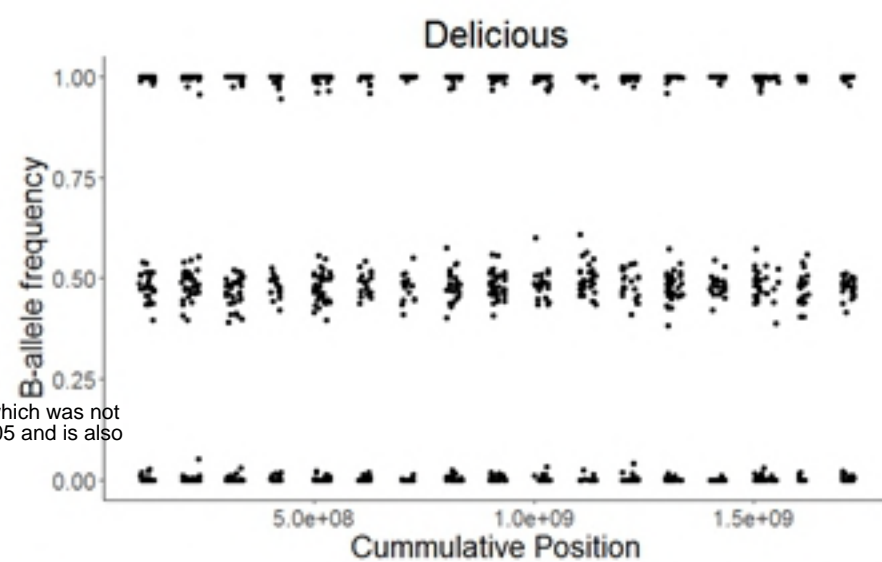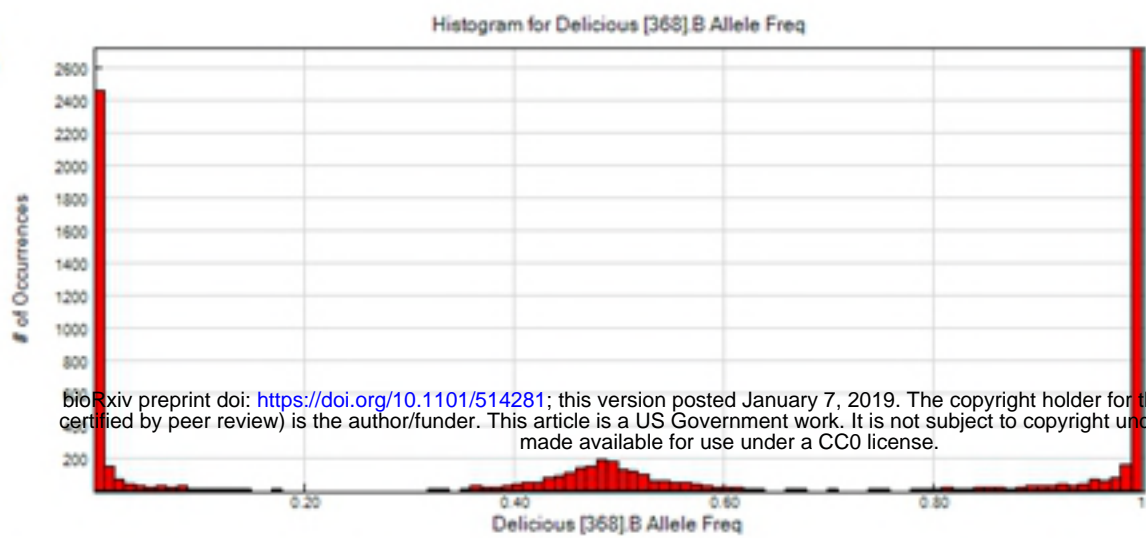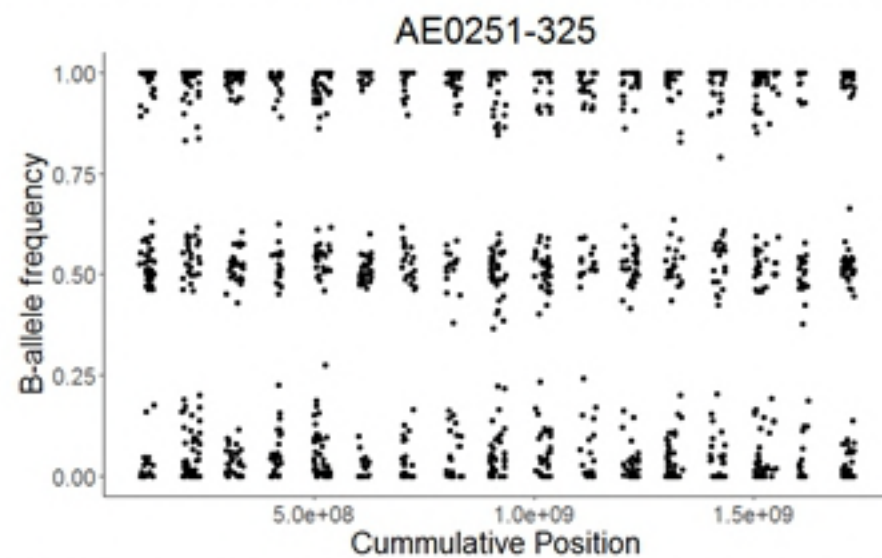
Figure 1

Figure 2