

RESEARCH ARTICLE

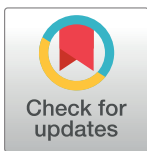
High-quality, genome-wide SNP genotypic data for pedigreed germplasm of the diploid outbreeding species apple, peach, and sweet cherry through a common workflow

Stijn Vanderzande¹*, Nicholas P. Howard^{2,3}, Lichun Cai⁴, Cassia Da Silva Linge⁵, Laima Antanaviciute⁵, Marco C. A. M. Bink^{6,7}, Johannes W. Kruijselbrink⁶, Nahla Bassil⁸, Ksenija Gasic⁵, Amy Iezzoni⁴, Eric Van de Weg⁹, Cameron Peace¹

1 Department of Horticulture, Washington State University, Pullman, WA, United States of America, **2** Department of Horticultural Science, University of Minnesota, St Paul, MN, United States of America, **3** Institute of Biology and Environmental Sciences, Carl von Ossietzky Universität, Oldenburg, Germany, **4** Department of Horticulture, Michigan State University, East Lansing, MI, United States of America, **5** Department of Plant and Environmental Sciences, Clemson University, Clemson, South Carolina, United States of America, **6** Biometris, Wageningen UR, Wageningen, The Netherlands, **7** Research & Technology Center, Hendrix Genetics, Boxmeer, The Netherlands, **8** USDA-ARS, National Clonal Germplasm Repository, Corvallis, OR, United States of America, **9** Plant Breeding, Wageningen UR, Wageningen, The Netherlands

* These authors contributed equally to this work.

* stijn.vanderzande@wsu.edu



OPEN ACCESS

Citation: Vanderzande S, Howard NP, Cai L, Da Silva Linge C, Antanaviciute L, Bink MCAM, et al. (2019) High-quality, genome-wide SNP genotypic data for pedigreed germplasm of the diploid outbreeding species apple, peach, and sweet cherry through a common workflow. *PLoS ONE* 14 (6): e0210928. <https://doi.org/10.1371/journal.pone.0210928>

Editor: David A. Lightfoot, College of Agricultural Sciences, UNITED STATES

Received: January 4, 2019

Accepted: April 19, 2019

Published: June 27, 2019

Copyright: This is an open access article, free of all copyright, and may be freely reproduced, distributed, transmitted, modified, built upon, or otherwise used by anyone for any lawful purpose. The work is made available under the [Creative Commons CC0](https://creativecommons.org/licenses/by/4.0/) public domain dedication.

Data Availability Statement: The private data might be obtained with the permission of the germplasm owners by contacting the crop's respective RosBREED team leader. Please contact Jim Luby for apple and Ksenija Gasic for peach. Jim Luby can be reached at lubyx001@umn.edu and Ksenija Gasic can be reached at kgasic@clemson.edu.

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

Funding: This project was co-funded by the USDA-NIFA-Specialty Crop Research Initiative projects, “RosBREED: Enabling marker-assisted breeding in Rosaceae” (2009-51181-05808), “RosBREED: Combining disease resistance with horticultural quality in new rosaceous cultivars” (2014-51181-22378), USDA NIFA Hatch projects 0211277 and 1014919, and the FruitBreedomics project No 265582: “Integrated approach for increasing breeding efficiency in fruit tree crops” that was co-funded by the EU seventh Framework Programme. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

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.

Introduction

A high-quality, mostly error-free genotypic data set is imperative to obtain reliable results in many downstream genetic analyses. The results of genetic analyses can be influenced by even low rates of genotyping errors [1]. For example, the size of genetic maps and order of markers therein are affected by errors in genotypic data [2–4]. Inaccurate genotypic data will also lower the power, accuracy, and resolution of linkage studies and increase the number of false marker-locus-trait associations [5–7]. The number of observed (double) recombinants is inflated by errors in genotypic data [8]. Incorrect calling of recombinations in turn leads to incorrect determination of haploblock limits and assignment of haplotypes [9]. Finally, incorrect genotype calls can lead to incorrect imputations of missing data or even the improper adjustment of correct data to ensure the data is consistent with Mendelian inheritance [10].

There are several reasons for the occurrence of errors in a genotypic data set. Incorrect information about a sample’s identity, e.g., due to mixing up or mislabeling samples, causes an individual to be matched with the wrong data [1]. In clonally propagated crops, mislabeling errors can easily spread when individuals that are not true-to-type are used as parents or as base plants to create new propagules. Available pedigree information for an individual can be incorrect, causing incorrect enforcement of allele assignments. In fruit cultivars, numerous pedigree records have been confirmed or updated with the help of genetic markers [11–23]. Biological reasons such as unexpected mutation, insertions or deletions in the DNA sequence containing markers, and gene conversion can lead to inconsistencies in genotype calls and propagate errors through the data set [1]. Technician errors can also introduce errors in a data set, such as when lab protocols are not applied correctly [24] or when multiple large data sets with disparate formats are integrated and edited. Finally, technological and software limitations and failures can also lead to the presence of errors [1].

SNPs have become the genetic marker of choice for many genetic analyses but, with their increased use and increasingly large numbers that can be generated, manual data curation has become more challenging. SNPs are ubiquitous within the genome and allow for simultaneous screening of many thousands of polymorphic loci via SNP arrays, Genotyping-By-Sequencing, or resequencing [25,26]. SNP arrays provide consistent information between individuals and have been developed for clonally propagated crops, such as the 8K apple array [27], 9K peach array [28], and 6K cherry array [29] developed by international teams led by RosBREED; the GrapeReSeq 18K *Vitis* array [30]; the 20K apple array developed by FruitBreedomics [31], all on the Illumina Infinium platform, and the strawberry 90K Axiom array [32], and the 480K apple array by FruitBreedomics on the Affymetrix axion platform [33]. Genotyping each individual relies on the automated scoring of thousands of SNPs. As thousands to millions of SNPs are being assessed on a large set of individuals, even a low error rate in SNP scoring can correspond to a high absolute number of errors. As the number of SNPs on an array increase, it becomes more time-consuming and less feasible to manually review all automated SNP calls to identify potential errors.

For SNP arrays, incorrect genotype assignment using automated SNP scoring software occurs when intensity plots deviate from expected patterns. Automated genotyping is based on the association of specific alleles to different fluorescent molecules, the detection of these fluorescent molecules, the clustering of individual-marker data points according to intensity ratios between the different fluorescent dyes across multiple individuals into distinct regions of a genotype-calling space, and the final assignment of these clusters to genotypes. Examples of deviations that are observed in the intensity plots are the presence of additional clusters or clusters that have shifted from their expected location in the intensity plot. The presence of additional clusters or shifted clusters can be attributed to additional regions that bind to the SNP's probe [34]. Sequence similarity of these regions with the intended target is caused by either local sequence repetition or presence of paralogous regions in the genome. The presence of these highly similar sequences can lead to multi-locus segregating SNP markers that cannot be adequately called. The calling of a single segregating locus might also be hampered by the background signal of targeted but non-segregating gene copies (ASSIST Reference Manual p17 [35]). The presence of one or more additional SNPs, insertions, or deletions in the probe-binding region can lead to reduced or loss of binding affinity for the SNP's probe and thereby to the presence of additional clusters, both of which can lead to incorrect genotype scoring of some SNPs [34].

No systematic workflow exists to efficiently detect and resolve all types of errors from a genotypic data set for pedigreed germplasm. Methods and software exist to tackle specific types of errors. For example, the ASSIST software was developed for use with Illumina Infinium arrays to identify which SNPs show robust results, which SNPs might have genotype calling errors due to alleles with reduced affinity or null alleles, and which SNPs are monomorphic or failed completely [36]. Another example is the aggregation of linked SNPs into a single genetic locus, called haploblock, which facilitates tracking the inheritance of alleles within a pedigree and subsequent identification of inheritance inconsistencies [37]. Despite the existence of these and other methods and software, an effective way to combine these methods has not been described.

Here we describe a curation workflow for high-resolution genetic marker data that identifies and resolves errors to obtain a robust set of genotypic data. The workflow maximizes the genotypic data obtained from high-throughput genome-scanning tools while minimizing the time needed to identify and remove errors. The workflow resulted from curation needs in the multi state and multi-crop USDA-SCRI project RosBREED [38–40] and the European project FruitBreedomics [41–43]. The workflow is demonstrated for three tree fruit crops, apple, peach, and sweet cherry, using the RosBREED germplasm sets [44]. The resulting genotypic data sets can be used by researchers to reconstruct pedigrees, establish quantitative genetic relationships, identify and validate quantitative trait loci (QTLs), and trace allele sources, leading to valuable practical and scientific genetic insights—with high confidence in the obtained results.

Material and methods

Plant material

The apple, peach, and sweet cherry collections used in this study, referred to as the 'Crop Reference Sets', were created to represent U.S. breeding germplasm [44] for the RosBREED project [38] (www.rosbreed.org) and consisted of 421, 426, and 269 individuals for apple, peach, and sweet cherry, respectively (S1–S3 Tables). Three apple breeding programs (Washington State University, the University of Minnesota, and Cornell University), three peach breeding programs (University of Arkansas, Clemson University, and Texas A&M University), and one

sweet cherry program (Washington State University) each contributed additional germplasm to complement the Crop Reference Sets and better represent their important breeding parents [44]. These additional ‘Breeding Pedigree Sets’ consisted of 176, 163, and 169 apple individuals, 117, 289, and 143, peach individuals, and 259 sweet cherry individuals, respectively. The sweet cherry Breeding Pedigree Set was later made publicly available and became part of the sweet cherry Crop Reference Set. Genotypic data of the other Breeding Pedigree Sets were included as part of the data curation, but individual identities of this private germplasm are not provided.

To reduce the trimming of pedigrees (as described under ‘Haploblock and haplotype generation’ below), the genotype calls of 18 additional apple individuals genotyped with the 20K SNP array in the FruitBreedomics project [43] or genotyped with the 8K SNP array at KU Leuven, Belgium (S1 Table) were added to the data set to complete genotypic data of key ancestors.

Initial parentage information

Initial parentage information was collected as part of the germplasm creation as described by Peace and co-workers (2014) [44]. For each breeding program, breeders provided pedigree records for their seedlings, selections, and released cultivars. Other pedigree records were based on historical records and available literature and were included for all progenitors, regardless of availability so that all progenitors terminated in founders (individuals with two unknown parents).

DNA extraction, SNP arrays, and iScan

DNA extraction was conducted for apple, peach, and sweet cherry as described by Chagné and co-workers (2012) [27], Verde and co-workers (2012) [28], and Peace and co-workers (2012) [29], respectively. Genomic DNA from each individual was purified using the E-Z 96 Tissue DNA Kit (Omega Bio-Tek, Inc., Norcross, GA, USA). DNA was quantitated with the QuantiT PicoGreen Assay (Invitrogen, Carlsbad, CA, USA), using the Victor multiplate reader (Perkin Elmer Inc., San Jose, CA, USA). DNA concentrations were adjusted to a minimum of 50 ng/μl, in 5 μl aliquots. For apple, DNA samples were run on the Illumina Infinium 8K apple SNP array [27] with iScans either at the Biotechnology Platform of the Agricultural Research Council (Pretoria, South Africa) or at the Research Technology Support Facility at Michigan State University (East Lansing, MI, USA), following the manufacturer’s protocol (Illumina Inc.). For peach and sweet cherry, DNA samples were run on the 9K peach SNP array [28] and 6K cherry SNP array [29], respectively, with an iScan at the Research Technology Support Facility at Michigan State University (East Lansing, MI, USA), following the manufacturer’s protocol (Illumina Inc.).

All three SNP arrays used in this study were developed through international collaboration led by the RosBREED project [27–29,38]. For each array, a set of individuals representing the worldwide breeding germplasm was chosen for next-generation sequencing to identify putative (*in silico*) SNPs throughout the genome. A subset of putative SNPs for each crop was first validated using the Illumina GoldenGate assay to fine-tune the filter parameters used in the final choice of SNPs for the array. In apple, 27 individuals were sequenced and their sequence was aligned to the ‘Golden Delicious’ apple reference genome v1.0 [27,45]. This genome version was described as a “high-quality draft” with a N50 of 16.7 kb. However, approximately 25% of the contigs had an “uncertain orientation” and genetic positions in linkage maps of 11% to 36% markers did not match with the markers’ physical positions [45,46]. Despite the genome’s shortcomings, it was successfully used to develop the apple 8K SNP array which

included 7758 putative SNPs [27]. The original evaluation of the apple 8K SNP array using 1616 apple individuals indicated that approximately 5554 (71%) of the putative SNPs were polymorphic [27]. In peach, whole-genome sequencing was obtained for 65 individuals and aligned to the ‘Lovell’ double haploid peach genome v1.0 to detect SNPs [28,47]. The peach genome v1.0 represented a high-quality whole genome shotgun chromosome-scale assembly with high contiguity (contig L50 214.2 kb), large portions of mapped sequences (96%) and high base accuracy (99.96%). However, few misassembly and orientation issues were detected and fixed in the peach genome v2.0 [48]. A total of putative 8144 SNPs were included on the peach 9K SNP array and its evaluation, carried out with 709 accessions from Europe and U.S., demonstrated that approximately 86.9% of the putative markers were polymorphic in the germplasm [28]. For cherry, a SNP detection panel of 16 sweet and 8 sour cherry accessions was chosen for whole genome, low-coverage resequencing [29]. Due to the lack of a reference genome for sweet cherry or tart cherry, the peach genome v1.0 [47] was used as a proxy for the cherry genome. The array included 4214 putative SNPs targeting the sweet cherry genome and 1482 putative SNPs targeting the sour cherry genome (both *avium* and *fruticosa* subgenomes) [29]. An evaluation of the cherry 6K SNP array using 269 sweet cherry individuals verified 1825 SNPs as polymorphic in sweet cherry breeding germplasm [29].

Initial genetic maps

For each crop, available genetic maps were used as a framework to determine the initial order of reliable SNPs. Reliable SNPs (obtained as described under ‘Subset of reliable SNP obtainment’ below) that were not present in available genetic maps were incorporated by comparing their physical positions to those of flanking SNPs that were present in the available genetic maps. For apple, information from an integrated map based on five full-sib families [20], the iGL map [49], and the “Golden Delicious” double haploid genome sequence v1.1 [50] were combined to generate the genetic map that was used in this study (‘Initial genetic map for apple’ in [S1 File](#)). The peach physical map was scaled to an approximate genetic map by using a conversion factor where every 1 Mb corresponded to 4 cM. For sweet cherry, genetic positions were determined by aligning and integrating the physical positions using peach genome v2.0 [48] with the sweet cherry ‘Regina’ × ‘Lapins’ SNP linkage map [21,51].

Workflow procedures

Throughout the workflow, several software packages were used. Below are described the main procedures used in the workflow, the associated software and parameter settings, and output files used. The order in which each functionality was used in the workflow is reported in Results section ‘Steps of the data curation workflow’.

Initial genotypic data obtainment (GenomeStudio). iScan output was converted to ‘AA’, ‘AB’, and ‘BB’ genotype calls for each SNP marker with the Genotyping module of GenomeStudio v2011.1 (Illumina Inc., San Diego, CA, USA) using a sample sheet (‘Sample sheet preparation’ in [S1 File](#)) to load sample intensities and a ‘Gen Call’ Threshold of 0.15 to assign samples to a genotype cluster. Genotypes in GenomeStudio were assigned as AA, AB, or BB where the A allele represents an A or T nucleotide and the B allele represents a C or G nucleotide depending on the assayed SNP.

Low-quality and non-diploid sample identification (GenomeStudio and R). iScan quality and ploidy levels were assessed using each sample’s B-allele frequencies calculated by GenomeStudio. This B-allele frequency represents the proportion of intensity observed for a sample that can be attributed to the B allele of the SNP. In GenomeStudio, the histogram of the B-allele frequencies observed across loci was plotted for each individual by using the ‘Histogram plot’

function of the ‘Full Data Table’ (‘B-allele frequency plots’ in [S1 File](#)). Samples were considered of good quality when a clear heterozygous peak was observed around 0.5 with almost no SNPs having a B-allele frequency between 0.125 and 0.375 and between 0.625 and 0.75. In contrast, samples of poor quality showed no clear heterozygous peak around 0.5 and had many SNPs with a B allele-frequency between 0.125 and 0.375, and between 0.625 and 0.75. Individuals that showed more than three peaks in the histogram were classified as polyploid. Individuals that showed a ‘shoulder’ on the AB peak were classified as putative aneuploids and were examined further in B-allele frequency plots according to Chagné and co-workers (2015) [52]. To create such a plot for an individual, the individual’s B-allele frequencies were exported from GenomeStudio for a subset of SNPs and plotted against the SNPs’ cumulative genetic positions (‘B-allele frequency plots’ in [S1 File](#)).

In the frequency plots according to Chagné and co-workers (2015) [52], ‘B-allele frequency’ values were expected to be in the same intervals as for the histogram plots. Diploid samples were considered of sufficient, intermediate, or poor quality when less than 0.3%, between 0.3% and 3%, or more than 3% of the subset of SNPs were observed between 0.125–0.375 and 0.625–0.875. For triploids, ‘B-allele frequency’ values were 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 expected for tetraploids. Aneuploids had a diploid pattern for most chromosomes and a haploid or polyploid pattern for others. Individuals classified as poor quality, polyploid, and aneuploid were excluded from further analyses.

Samples were excluded from various input files and from the genotype clustering in GenomeStudio by choosing them in the ‘Samples Table’ and then choosing the ‘Exclude Selected Samples’. SNPs were then re-clustered by choosing the ‘Cluster All SNPs’ of the ‘Analysis’ section. All statistics were updated when prompted.

Subset of reliable SNP obtainment (ASSIsT). The ‘Final report’, ‘DNA report’, ‘pedigree’, and ‘map’ input files were created as described in the ASSIsT Reference Manual [35] and the ‘ASSIsT input files’ section of [S1 File](#). All input files were loaded into ASSIsT v1.01 [36] and parameters were set depending on the ‘Population type’ ([S4B Table](#)). Eight marker classes distinguished by ASSIsT, re-grouped into the following five categories:

- Robust SNPs: Less than 5% No Call Rate and all three possible clusters (AA, AB, and BB) present in the germplasm set. In ASSIsT, these SNPs were classified as ‘Robust’, ‘OneHomozygRare_HWE’, ‘OneHomozygRare_NotHWE’, and ‘DistortedAndUnexSegreg’
- Two cluster SNPs: Less than 5% No Call Rate and one of the homozygous clusters (AA or BB) absent in the germplasm set. In ASSIsT, these SNPs were classified as ‘ShiftHomo’
- Null-allele SNPs: Probable presence of a null allele, classified as ‘NullAllele-Failed’ in ASSIsT. Null alleles are alleles that do not lead to an observable signal and either reduce the overall signal intensity of a sample (when the second allele is detectable) or lead to no signal at all (when the sample is homozygous for the null allele)
- Monomorphic SNPs: No polymorphism, as in ASSIsT
- Failed SNPs: Having more than 50% No Call Rate, poor clustering, or low intensity, as in ASSIsT

Results of SNP performance in ASSIsT were exported to the ‘Summary’ and ‘Custom SNP information table’. Genotype calls were saved in ‘Custom gtypes’ to be used in the R-script that checked pedigree records (described below in ‘Pedigree records verification’). PLINK input files were generated to check for unknown duplicates within the data (described below in ‘Duplicate individuals detection’) and FQ_DataPrepper input files were created to easily

generate FlexQTL input files using FQDataPrepper (described below in ‘Genotyping error detection and adjustment’). Genotype calls for the ‘Robust SNPs’ category were automatically reported in ASSiST output files whereas other categories were considered to contain failed SNPs and thus their genotype calls were not automatically reported. To include genotype calls of the ‘Two cluster SNPs’, genotype calls of such SNPs were extracted from GenomeStudio and added to the data files manually.

Duplicate individuals detection (GenomeStudio and Plink). Genotypic data of known mutants and duplicates were compared to ensure their genotypic data were matching using the ‘Reproducibility and Heritability’ report of GenomeStudio (Analysis>Reports>Reproducibility and Heritability Report>with Calculating Errors). The data set was also screened for individuals with (unknown) identical genotypic data using Plink 1.9 [53] (<https://www.cog-genomics.org/plink2>). Plink 1.9 was run using input files generated with ASSiST to calculate the proportion of identity-by-descent (IBD) between each pair of individuals (‘Plink analysis’ in [S1 File](#)). Pairs of individuals with an IBD proportion higher than 97% were considered to be duplicates because at this stage all known duplicates shared an IBD proportion of at least 97%. If individuals were true duplicates, only one was kept in the data set. If pedigree records differed between duplicate individuals, pedigree records were used to identify trueness-to-type as described below. True-to-type individuals were kept in the data set and individuals that were not true-to-type were targeted for DNA re-sampling. Where two unselected seedlings from the same family were identified as duplicates, they were both targeted for re-sampling as it was unclear which of the two was true-to-type.

Pedigree records verification (GenomeStudio, Cervus, and R). Verification of pedigree records was performed by counting the Mendelian-inconsistent errors between an individual and (each of) its recorded parent(s) where genotypic data was available. These errors were genotypic data inconsistent with Mendel’s first law, i.e., alleles present in offspring but not present in either parent. First, parent-child (PC) errors between an individual and a single parent were defined as genotype calls where none of the parental alleles were present in the offspring. For example, the recorded offspring might be ‘BB’, ‘B null’, or ‘null null’ while the recorded parent was ‘AA’. In this example, neither the ‘B’ allele nor the ‘null’ alleles were present in the parent. Secondly, when both parents were known and confirmed, the combination of the two parents’ SNP data were compared to the offspring’s SNP data to identify parent-parent-child (PPC) errors. PPC errors were defined as genotype calls where at least one allele of the offspring was not present in any of its recorded parents. For example, in the case of an ‘AA’ x ‘AA’ -> ‘AB’ triplet, no PC error would be observed when checking each parent individually, as both parents could have contributed the ‘A’ allele to the offspring. However, combination of the two parents would create a PPC error as neither parent could have contributed the ‘B’ allele observed in the offspring.

Three approaches to count Mendelian-inconsistent errors were compared. In GenomeStudio, a ‘Reproducibility and Heritability’ (Analysis>Reports>Reproducibility and Heritability Report>with Calculating Errors) was generated to obtain the number of PC and PPC errors. Mendelian-inconsistent errors were calculated in the software Cervus [54] using default parameter settings. Third, an ad hoc R-script ([S3 File](#)) was used to check and identify PC and PPC relationships as described in the ‘R script for PC and PPC relationships’ section of [S1 File](#).

When an individual’s supposed parent was not genotyped but the supposed grandparents were genotyped, the grandparents-grandchild relationship was tested with the AB+AA-AA test in Excel using the template provided by van de Weg and co-workers (2018) (‘AB+AA-AA test’ in [S1 File](#)) [23].

A threshold was determined for the proportion of PC errors to confirm or reject PC relationships using incompletely curated marker data. PC errors were counted for 1000 pairs of two

random individuals in the data set that did not have a (known) PC relationship and for all pairs of individuals that had a known PC relationship. A separation was observed between the resulting distributions of PC errors for the two sets of individuals and a midway point between both distributions was used as threshold to reject the parentage of an individual. Similarly, a threshold was determined to accept or reject the combination of two parents; observed PPC errors were counted for previously confirmed PPC relationships and a threshold set as 110% of the highest number observed PPC errors among these known relationships.

In case of missing or erroneous parent information, efforts were made to identify the missing parent and, if not possible, to identify sets of possible grandparents. Hereto, all available selected material was examined (ancestors, direct parents, and breeding selections). In apple and peach, the ad hoc R-script ([S3 File](#)) was used to find PC and PPC relationships ('R script for PC and PPC relationships' in [S1 File](#)). In cherry, the software Cervus [54] was used to count these errors and determine possible parents using the default parameter settings. When no second possible parent was found in the data set, possible grandparents were identified in Excel using the template provided by van de Weg and co-workers (2018) [23]. Historic records (e.g., location and time of origin) of possible grandparents were checked to ensure feasibility. Furthermore, deduced grandparent-grandchild relationships were only kept if they did not lead to a large number of reported errors during the rest of the workflow.

Pedigree information was then updated in various input files and in GenomeStudio (Analysis>Edit Parental Relationships; then choosing individual and correct parents from drop-down menu) for further analyses. All statistics in GenomeStudio were updated when prompted.

Genotyping error detection and adjustment (GenomeStudio, FlexQTL, and Visual FlexQTL). Genotyping errors were divided in two classes: Mendelian-inconsistent errors and Mendelian-consistent errors [10]. Unlike Mendelian-inconsistent errors, Mendelian-consistent errors are errors that do not infringe upon Mendel's first law: a child's false allele call is present in one of the parents, but results in problematic co-segregation patterns that show unexpected double recombination between markers with successive genetic/physical positions. These double recombinations might be due to issues in ploidy, genotype calling, marker ordering, or phasing or, occasionally, gene conversion [10] ([S4 File](#)).

For individuals with verified pedigree relationships, remaining Mendelian-inconsistent errors were detected using GenomeStudio and FlexQTL v0.99130. In GenomeStudio, the 'SNP Table' was filtered for SNPs with Mendelian-inconsistent errors, the 'Error Table' was used to identify individuals with Mendelian-inconsistent errors, and the 'SNP Graph' was used to examine the reported errors. FlexQTL input files were prepared using FlexQTL DataPrepper v1.0.0.4 (<https://www.wur.nl/en/show/FlexQTL.htm>) and the generated data file was further adjusted to ensure all individuals had either both parents specified or none ('FlexQTL input files' in [S1 File](#)). FlexQTL was used to check for Mendelian-inconsistent errors (parameter settings in [S4C Table](#)). Briefly, FlexQTL was run through using an early stop ('pedimapV' parameter set to '2'; to stop after checking the data for inconsistencies) and allowing for segregation distortion ('MSegDelta' parameter set to 1). This analysis summarized for each marker and each individual how many Mendelian-inconsistent errors were observed in the 'mconsistency.csv' file.

Mendelian-consistent errors were detected by examining double recombinations detected over small regions (<10 cM) as reported by FlexQTL and Visual FlexQTL. Parameter settings of FlexQTL to check for double-recombinations were the same as for Mendelian-inconsistent errors above ([S4C Table](#)). The FlexQTL output file named 'DoubleRecomb.csv' listed all singletons (single markers involved in a double recombination) in the data set. Visual FlexQTL instead identifies all double recombinations (including singletons) that occur within a given

genetic distance. The default for this distance was 10 cM and could be changed under ‘Tools>Calculate>(Re-)Compute recombination sequences’. The report on double recombinations was created through ‘Tools>Export>Export recombination sequence file’ which provided an output file called ‘DoubleRecombinations.csv’.

Genotype calls of SNPs with Mendelian-inconsistent errors or SNPs involved in detected double recombinations were further examined in GenomeStudio using the ‘SNP Graph’. Where incorrect cluster identification was detected, clusters were manually called using the ‘SNP Graph’ (‘manual SNP calling’ section in [S1 File](#)) and FlexQTL was run again to ensure errors were resolved. If genotype calls could not be accurately made, the SNP was considered to have failed and removed from the data set.

Identification of Mendelian-inconsistent and Mendelian-consistent errors were also performed at the haplotype level, conducted as described above at the single SNP level. Where an unidentified error in SNP genotype scoring was detected, the corresponding SNP genotype calls were adjusted. If the calling error occurred in a single or few individuals, haplotypes were manually adjusted to reflect the change in SNP allele. In the rare event that a large group of individuals had their SNP genotype calls adjusted, the corresponding haplotypes were re-determined using PediHaplotyper [37]. Where Mendelian-inconsistent errors were due to missing SNP alleles, the individual was compared to its parent and offspring to determine the correct haplotype. For example, if an individual had a SNP haplotype of ‘A-?-B-A’ and the haplotype was not present in either parent, but a parent had a haplotype of ‘A-A-B-A’ and 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 ‘A-B-B-A’ were present in the parent, information of flanking, linked haplotypes were checked to assess if the offspring’s haplotype could be determined by minimizing the number of recombinations. Where inconsistencies in selected material were suspected to be due to a recombination in an ungenotyped progenitor, the haploblock was split in two at the suspected recombination site to avoid inconsistency in downstream genetics analyses of recombination in selected material. The haplotypes for those two new haploblocks were determined again using PediHaplotyper.

Map error detection and adjustment (FlexQTL, Visual FlexQTL, and Microsoft Excel). Where double recombinations were observed and these recombinations were not due to incorrect genotype scoring, a graphical genotyping approach was used to examine and possibly adjust SNP order in the genetic map [55]. Graphical genotyping plots were created starting from the ‘SIP_Population.csv’ output file of FlexQTL ([S4 File](#)). FlexQTL was run again to ensure the errors were resolved and only if the adjustment of the SNP order did not lead to new double recombinations, a change in order was accepted. SNPs were removed from the data set if they had unexpectedly high incidences of double recombinations that could not be resolved by repositioning the SNPs in the map. Additionally, where a SNP mapped to multiple locations in different families, the SNP was removed from the data set.

Haploblock and haplotype determination (FlexQTL, Visual FlexQTL, and PediHaplotyper). Haploblocks were defined as regions in which no recombination was observed for selected material. Pedigree data was trimmed to ensure accurate phasing and dummy individuals were introduced to ensure all selected material is considered during haploblock border determination (‘Haploblock and haplotype determination’ in [S1 File](#)). The data was phased using FlexQTL (parameter settings in [S4D Table](#)) and haploblock borders were defined with Visual FlexQTL (‘Haploblock and haplotype determination’ in [S1 File](#)). After removing the dummy individuals introduced for haploblock border determination, a second run of phasing with FlexQTL (parameter settings in [S4D Table](#)) provided the remaining input files for PediHaplotyper (‘Haploblock and haplotype determination’ in [S1 File](#)).

The PediHaplotyper package [37] was run in R and the resulting output files ('Hapblock and haplotype determination' in [S1 File](#)) were used as input files for FlexQTL for further data curation of the haplotyped data sets (resolving both Mendelian-inconsistent and Mendelian-consistent errors as described under '*Genotyping error detection and adjustment*').

SNP classification. A SNP classifications system was established to track clustering issues and minimize future curation of new data. SNPs that passed the filter criteria from ASSIsT and that were included in the final data set were classified into four types: type 1 SNPs had no or less than 5% call editing during the curation process and no additional genotype clusters were present; type 2 SNPs had an incorrect automated cluster identification of one of the genotype clusters (e.g., 'AA' cluster called as 'AB'), showed no additional clusters, and could easily be corrected; type 3 SNPs showed additional clusters because of alleles with differential intensity signals but individuals could easily be called correctly; and type 4 SNPs had null alleles but individuals with null alleles could be distinguished easily from true homozygous individuals. Type 5 SNPs could be accurately called but their genetic or physical position could not be determined accurately and were not included in the map and final data set. Type 6 SNPs were monomorphic across all individuals. Type 7 SNPs were those considered as 'Failed' by ASSIsT or were removed during the workflow because their genotype calls could not be manually resolved.

Workflow creation and implementation

A workflow was constructed by identifying necessary steps of data curation and ordering them in such a way that the amount of time needed for data curation is minimized at each step. Thus, errors addressed first were those relatively easy to identify and resolve and otherwise expected to cause problems at multiple steps. The workflow was an outcome of efforts in RosBREED and FruitBreedomics on data curation in apple, peach, and cherry. Statistics at each step of curation were determined from implementing this workflow on the RosBREED germplasm described in the 'Plant Material' section above.

Results

Steps of the data curation workflow

Initial error-detection resulted in a list of possible causes for each type of detected errors ([Table 1](#)). This list identified which issues had to be resolved first and as such resulted in the workflow described below ([Fig 1, S4 File](#)). The workflow developed had three main parts, each with multiple steps. The first main part ensures that genetic principles can be applied, the second main part applies these principles on a single marker level, and the last main part applies these principles at the haploblock level. The proposed steps within each main part are described below, as conducted for apple, peach, and sweet cherry.

1. Ensuring inheritance principles can be applied. After creating an initial data set of genotypic data set in GenomeStudio, a first set of analyses was performed. Because genotypic errors are identified based on principles of inheritance in diploids, individuals and markers that do not follow these principles had to be removed first ([Fig 1](#)). When doing so, individuals with unexpected intensity patterns had to be removed first ([Fig 1](#)) as they were influencing the clustering of all individuals in the germplasm. Individuals with poor quality DNA were usually poorly genotyped, resulting in many data inconsistencies. Additionally, polyploids (individuals having one or more additional full chromosome sets) and aneuploids (individuals having an irregular number of copies for one or more chromosomes) were expected to have intensity ratios for heterozygous loci that differed from diploid individuals. Removal of individuals with

Table 1. Errors observed during the curation process and their possible causes. Causes that should be (mostly) already resolved by the stage a researcher would start checking for specific errors are in parentheses and italicized.

Error	Cause	Solution
Low call rate and impossible cluster identification	Probe binding issues	Remove SNP from data set
Unexpected B-allele frequencies	<i>(Probe binding issues)</i>	<i>(Remove SNP from data set)</i>
	Unexpected ploidy	Remove sample from data set
	Low sample quality	Remove sample from data set
High number P(P)C errors	<i>(Probe binding issues)</i>	<i>(Remove SNP from data set)</i>
	<i>(Low sample quality)</i>	<i>(Remove sample from data set)</i>
	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	<i>(Probe binding issues)</i>	<i>(Remove SNP from data set)</i>
	<i>(Low sample quality)</i>	<i>(Remove sample from data set)</i>
	<i>(Incorrect pedigree)</i>	<i>(Adjust pedigree record)</i>
	Incorrect clustering	Manually determine genotype clusters
	Incorrect genotype call(s) not due to cluster issues	Adjust genotype call(s)
High number double recombinations	<i>(Probe binding issues)</i>	<i>(Remove SNP from data set)</i>
	<i>(Low sample quality)</i>	<i>(Remove sample from data set)</i>
	<i>(Incorrect pedigree)</i>	<i>(Adjust pedigree record)</i>
	<i>(Unexpected ploidy)</i>	<i>(Remove sample from data set)</i>
	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	<i>(Probe binding issues)</i>	<i>(Remove SNP from data set)</i>
	<i>(Low sample quality)</i>	<i>(Remove sample from data set)</i>
	<i>(Incorrect pedigree)</i>	<i>(Adjust pedigree record)</i>
	<i>(Incorrect clustering)</i>	<i>(Manually determine genotype clusters)</i>
	<i>Nearby double recombination*</i>	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	<i>(Probe binding issues)</i>	<i>(Remove SNP from data set)</i>
	<i>(Low sample quality)</i>	<i>(Remove sample from data set)</i>
	<i>(Incorrect pedigree)</i>	<i>(Adjust pedigree record)</i>
	<i>(Incorrect clustering)</i>	<i>(Manually determine genotype clusters)</i>
	<i>(Incorrect marker position in map)</i>	<i>(Adjust marker position or remove marker if it cannot be accurately mapped)</i>
	<i>(Incorrect genotype call(s) not due to cluster issues)</i>	<i>(Adjust genotype call(s))</i>
	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

<https://doi.org/10.1371/journal.pone.0210928.t001>

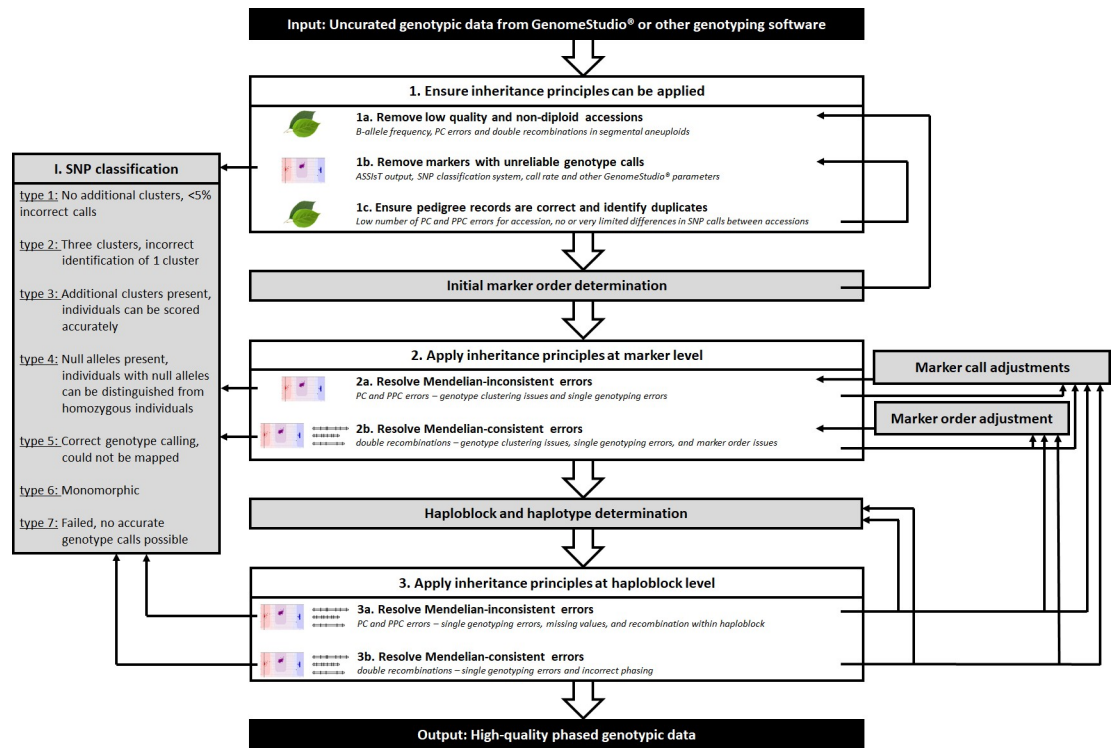


Fig 1. Steps of the high-resolution genotypic data curation workflow to ensure a quick and efficient curation process. Steps that identify errors are shown in white boxes; procedures needed for detecting, keeping track of, and resolving errors but do not identify errors directly are in grey boxes. After obtaining a first set of genotypic data, initial steps ensure that inheritance principles can be readily applied by removing individuals and markers that do not follow these principles and by ensuring pedigree records are correct. In the next set of steps, inheritance principles are applied at the individual marker level. In the final set of steps, these principles are applied at the haploblock level. Output used to detect and resolve observed errors at each step are given in italics. The leaf symbol indicates errors at the level of individual; the intensity plots symbol indicates errors at the level of SNP scoring; the genetic map symbol indicates errors at the level of genetically linked markers and phased alleles. When applying inheritance principles in parts 2 and 3, alleles that do not occur in an individual's parents ('Mendelian-inconsistent errors') are first resolved before addressing remaining genotyping errors ('Mendelian-consistent errors'). Several procedures, such as marker call adjustments and map order adjustments, are performed throughout the steps of the workflow to resolve errors detected. Each time after performing these common procedures, specific steps of the workflow must be repeated, forming an iterative process that ends when all errors are resolved.

<https://doi.org/10.1371/journal.pone.0210928.g001>

poor DNA quality and suspected polyploids and aneuploids was observed to improve genotype cluster definitions and thereby the genotype calling of remaining individuals.

Once individuals with ploidy and sample quality issues were removed, a set of well performing markers had to be obtained (Fig 1). Markers with unreliable scoring were observed to lead to many inconsistencies in subsequent steps. Thus, their early removal would ensure that a relatively low number of inconsistencies remained in the data set, expected to greatly reduce the observed inconsistencies and time needed for further steps.

Identifying and correcting incorrect PC and PPC relationships was a prerequisite to using pedigree information for the identification of marker calling errors in each data set. Imposing principles of inheritance on actually unrelated individuals led to many false errors at the marker and map level. Conversely, identifying thus far unknown PC and PPC relations helped to identify errors at the marker and map level elsewhere in the data set and was expected to improve the power of downstream QTL analyses. Thus, recorded pedigree information needed to be validated and previously unknown pedigree relationships deduced before curating individual marker calls and marker order errors (Fig 1). Duplicate individuals were also detected

at this stage as they could help resolve sampling errors and reduce the number of individuals needing detailed error-checking.

2. Applying inheritance principles at the marker level. When Mendelian-inconsistent errors were present, at least one allele was incorrect. This issue had to be resolved before the (corrected) allele could be phased with the alleles of flanking markers. Otherwise, even the other allele, which might have been correct, could have been incorrectly phased with the alleles of flanking markers, causing additional observed but false recombinations. Thus, to minimize the time required to resolve Mendelian-consistent errors by investigating many supposed double recombinations, Mendelian-inconsistent errors had to be addressed first.

Markers with a high number of errors were investigated before markers with a relatively low number of errors among progenitors. Then, markers with a low number of errors for seedlings were investigated as they were expected to have the least effect on the remaining data set.

Any supposed double recombinations that occurred at the same region in multiple individuals had to be resolved first as they were very unlikely, could be due to a single error, and could influence a large set of individuals. Next, suspicious double recombinations that occurred over multiple loci in ancestors had to be checked, followed by singletons in ancestors. Finally, singletons in seedlings were checked, but they were expected to be the least harmful when incorrect because of little to no effect on the remaining data set.

When no genotype calling or map errors were detected, phasing errors were investigated by checking the phasing of individuals that shared the parent whose homolog was observed to have a double recombination. In the rare case that incorrect phasing by FlexQTL led to a double recombination in multiple individuals of a single family or parent, it was always caused by one or two individuals in which the position of (a single) recombination was incorrectly determined. In those cases, individual(s) for which the SNP was involved in a single recombination had their genotype set to missing. This adjustment led to correct phasing of all other individuals and removal of reported double recombinations. Double recombinations that were observed in a single individual and that were not due to incorrect genotype clustering or incorrect map positions were accepted as the result of true double recombination events.

3. Applying inheritance principles at the haploblock level. Haploblock and haplotype determination was based on correctly identifying recombinations through correct phasing across generations and combining individual SNP alleles into haplotypes. Thus, any remaining errors at the SNP level or map level were expected to lead to errors in haploblock and haplotype determination. Therefore, all observed inconsistencies at the individual SNP level had to be resolved before inconsistencies were detected at the haploblock level. The genetic principles applied throughout the workflow are expected to also hold up at the haploblock level and therefore haplotypes had to be checked for Mendelian-consistent errors and Mendelian-inconsistent errors.

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

1a. Removing samples: Non-diploid individuals and low-quality samples. In apple, the 'B allele frequency' plot of 744 of the diploid individuals (80.7%) was very close to that expected for diploid individuals (Fig 2A; S1 Table) and results of these diploid individuals were considered to be of good quality. Another 71 individuals (7.7%) showed some variation from the expected B allele frequency, especially for homozygous SNPs, but the three genotypes could be easily distinguished (Fig 2B; S1 Table) and their results quality was considered to be intermediate. Finally, 107 (11.6%) had 'B allele frequency' plots that showed a wide variation around the

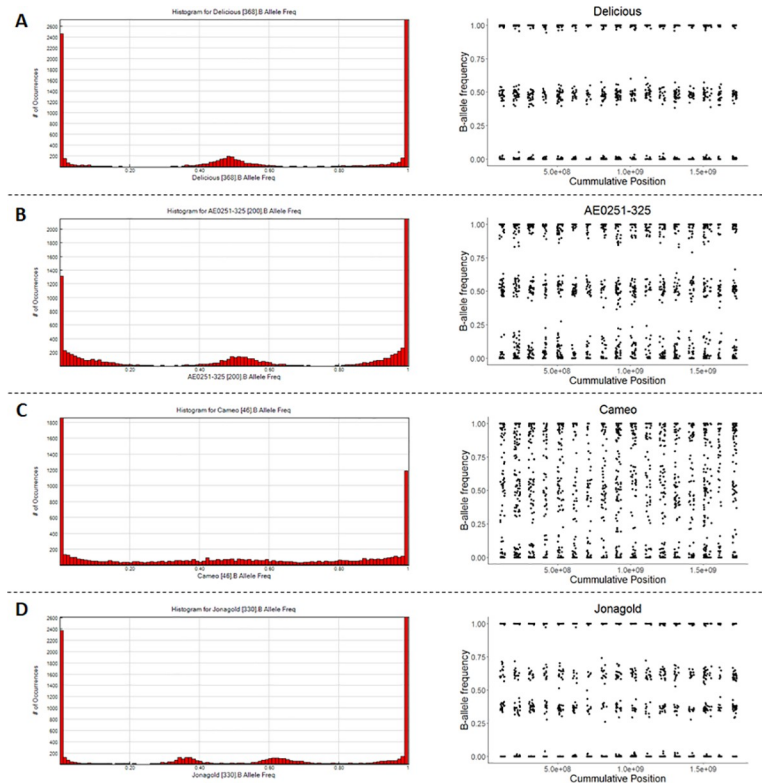


Fig 2. Histograms of B-allele frequency (left) and B-allele frequency for each SNP plotted against its genomic position (right). Such histograms were used to assess a sample’s genotyping quality and ploidy. Examples shown are of a sample with good quality genotype calls (panel A), with intermediate quality of genotype calls (B), with bad quality of genotype calls (C), and that is triploid (D).

<https://doi.org/10.1371/journal.pone.0210928.g002>

expected frequency (Fig 2C; S1 Table) and their results quality was considered to be bad. No individuals with bad quality results were found for peach or sweet cherry.

For apple, most individuals with poor quality results had their DNA extracts transported outside the U.S. for genotyping and the poor results were suspected to be caused by a reduction in DNA quality due to the delay in clearing customs, while only nine individuals with poor quality were from those genotyped in the U.S. The call rate in GenomeStudio differed between the individuals that had good, intermediate, or bad quality, with the call rate dropping as the level of quality lowered (S1 Fig).

For apple, five triploid individuals were identified (S1 Table). One was the known triploid cultivar ‘Jonagold’ while the others were unselected seedlings (S1 Table; S2A Fig). Two other unselected seedlings had their B-allele frequencies divided over 5 clusters of the GenomeStudio plot, which indicated they could be tetraploid or a mixture of two samples (S1 Table; S2B Fig). No aneuploids were detected in the apple germplasm. However, one individual from the Crop Reference Set, ‘AE213-200’ and one individual of a Breeding Pedigree Set were identified as segmental aneuploids (missing one copy of a large chromosomal segment). They were undetectable in the B-allele frequency analysis and instead identified by a relatively large number of PC errors and double recombinations observed for only that chromosomal segment. No polyploids, aneuploids, or segmental aneuploids were detected in peach and sweet cherry.

The final number of individuals used in the rest of the workflow was 835, 621, and 528 for apple, peach, and sweet cherry, respectively, consisting of 139, 48, and 56 direct parents of full-

Table 2. Summary of SNP classification by ASSIsT for apple, peach, and sweet cherry. SNP classifications are grouped in retained and discarded SNPs.

SNP classification	Apple	Peach	Sweet Cherry
Retained SNPs			
<i>Robust SNPs</i>			
Robust	1434	743	373
OneHomozygRare_HWE	357	62	109
OneHomozygRare_NotHWE	366	188	161
DistortedAndUnexSegreg	1362	3696	555
<i>Two cluster SNPs</i>			
ShiftedHomo	914	1409	529
<i>Total</i>	4433	6098	1727
Discarded SNPs			
NullAllele-Failed	52	145	43
Monomorphic	708	1057	3478
Failed	2565	844	448
<i>Total</i>	3325	2046	3969
Total	7758	8144	5696

<https://doi.org/10.1371/journal.pone.0210928.t002>

sib families, ancestors, and cultivars, 76, 24, and 9 selections and 620, 548, and 463 unselected seedlings over 45, 26, and 41 families of 4–62 full-sibs, respectively (S1–S3 Tables).

1b. Obtaining a set of reliable SNPs. A subset of SNPs with reliable genotyping scores was obtained using ASSIsT (Tables 2 and 3). Although discarded by ASSIsT, SNPs from the ‘Two cluster SNPs’ category were retained as many of them were considered to contain useful

Table 3. Summary of SNPs retained and discarded for apple, peach and sweet cherry during the steps of the workflow.

SNP curation step	SNPs discarded	SNPs retained
<i>Apple</i>		
1b. Set of reliable SNPs with ASSIsT	4252	4536
2a. Mendelian-inconsistent error detection at SNP level	319	4217
2b. Mendelian-consistent error detection at SNP level	329	3888
*Removed SNPs due to mapping issues	15	
*Removed SNPs due to genotyping issues	314	
3. Error detection at haplotype level	33	3855
<i>Peach</i>		
1b. Set of reliable SNPs with ASSIsT	2046	6098
2a. Mendelian-inconsistent error detection at SNP level	231	5867
2b. Mendelian-consistent error detection at SNP level	1862	4005
*Removed SNPs due to mapping issues	156	
*Removed SNPs due to genotyping issues	1706	
3. Error detection at haplotype level	-	4005
<i>Sweet Cherry</i>		
1b. Set of reliable SNPs with ASSIsT	3969	1727
2a. Mendelian-inconsistent error detection at SNP level	47	1680
2b. Mendelian-consistent error detection at SNP level	63	1617
*Removed SNPs due to mapping issues	63	
*Removed SNPs due to genotyping issues	0	
3. Error detection at haplotype level	-	1617

<https://doi.org/10.1371/journal.pone.0210928.t003>

information. A total of 4536 (58%), 6098 (75%), and 1727 (30%) of the SNPs on the apple, peach, and cherry arrays, respectively, were maintained after filtering. Subsequent steps of the workflow reduced the number of SNPs in the final data set further to 3855, 4005, and 1617 for apple, peach, and sweet cherry, respectively. Thus 85%, 66%, and 91% of the SNPs retained after using ASSiST for apple, peach, and sweet cherry, respectively, resulted in high-quality data.

1c. Correcting pedigree information and identifying duplicates. The number of PC errors in apple between two randomly paired individuals without PC relationship averaged 195, with a minimum of 17 (comparison between two full-sibs) and 99% of these comparisons had more than 40 errors. In contrast, average and maximum number of PC errors between two related individuals with a known PC relationship was 2 and 17, respectively, and 99% of these comparisons had less than 10 PC errors. The threshold to reject a PC relationship was set at 23 errors, which roughly corresponded to 0.5% of total markers. For 103, 66, and 22 individuals, one recorded parent was incorrect in apple, peach, and sweet cherry respectively, and for 36, 14, and zero individuals, both recorded parents were incorrect. For 106, 1, and 19 of these individuals in apple, peach, and sweet cherry, one or both of the true parent(s) was found within the germplasm set. The final number of generations spanned by the corrected pedigrees was eight, nine, and six for apple, peach and sweet cherry, respectively.

2a. Finding Mendelian-inconsistent errors at the SNP level. FlexQTL summarized the number of Mendelian-inconsistent errors for each marker and each individual. In GenomeStudio, the 'SNP Table' would summarize the number of P(P)C errors for each SNP and a separate 'Error Table' had to be consulted to determine which individuals were involved in these errors. FlexQTL mostly reported the error under the parent, the R-script reported the error under the offspring, and the 'Error Table' of GenomeStudio reported the genotypes of both parent(s) and offspring. As a consequence, errors between a single parent and multiple of its offspring would be reported as one erroneous (parental) genotype in FlexQTL whereas GenomeStudio reported the error for each offspring. However, FlexQTL did identify errors between grandparents and grandchildren when the missing parental genotype could be imputed.

FlexQTL detected 1209, 2230, and 686 Mendelian-inconsistent errors distributed over 541, 760, and 47 SNPs in apple, peach, and sweet cherry respectively. In apple, GenomeStudio detected 10,201 PC errors and PPC errors over 2303 SNPs. Although GenomeStudio identified which pairs of individuals led to these errors, some of the detected Mendelian-inconsistent errors did not occur in the data set due to differences in genotype scoring between ASSiST and GenomeStudio. Mendelian-inconsistent errors could not always be resolved. A total of 319, 231, and 47 SNPs with Mendelian-inconsistent errors (59%, 30%, and 100% of SNPs with errors as reported by FlexQTL, respectively) were removed in the apple, peach, and sweet cherry data sets, respectively, because genotypes for these SNPs could not be accurately determined (Table 3). Before removal of these Mendelian-inconsistent errors, 41,717, 29,009, and 2505 double recombinations involving a single marker were detected in FlexQTL in apple, peach, and sweet cherry, respectively, through the 'DoubleRecomb.csv' file, whereas only 6177, 4905, and 1739, respectively, of these recombinations were observed after removal of all Mendelian-inconsistent errors.

2b. Identifying Mendelian-consistent errors at the SNP level. Most double recombinations that occurred in the same genomic region in many individuals could be resolved by adjusting incorrect marker calls. A total of 648, zero, and 209 markers in apple, peach, and sweet cherry, respectively, had one or more of their genotype calls adjusted to resolve double recombinations. Most other double recombinations that occurred in multiple families could be resolved by repositioning the marker in the genetic map using a graphical genotyping approach. In total, 115, zero, and zero SNPs were moved from their original position in the

map to resolve double recombinations for apple, peach, and sweet cherry, respectively. Many recombination events that occurred in a single or few individuals over a single marker were resolved by first resolving the double recombinations that occurred in many individuals. Most of the remaining double recombinations were solved by either changing single incorrect genotype call or adjusting marker order in the map. Only a few phasing issues were observed where (almost) all offspring of a founder showed a double recombination that could be resolved by adjusting the phase of the alleles in that founder. A total of 15, 156, and 63 markers were discarded for apple, peach, and sweet cherry, respectively, because they led to unresolvable map issues (Table 3). An additional 314, 1706, and zero markers were removed while resolving Mendelian-consistent errors in apple, peach and sweet cherry, respectively, because no accurate genotyping could be obtained (Table 3). The total number of remaining reported singletons was 68, 47, 51 for apple, peach, and sweet cherry, respectively, and these were considered to be true double recombinations.

During data curation, genetic maps were generated for each crop (S5–S7 Tables) by adding new SNPs to existing maps, by converting physical positions into genetic positions, and/or by updating initial genetic positions to minimize the number of double recombinations. For apple, 885 SNPs were added and 658 previously-mapped SNPs were removed as they did not perform well in our wider germplasm. Addition of SNPs at the chromosome ends enlarged the original map by 7 cM. The resulting apple map was 1179 cM long with chromosome lengths ranging from 57.6 cM (linkage group (LG) 6) to 103.6 cM (LG 15). The number of SNPs on each LG ranged from 167 SNPs on LG 6 to 359 SNPs on LG 2. The genetic map of peach was 893.2 cM long; LG 5 was the shortest (72.9 cM) and LG 1 was the longest (190.2 cM). The number of SNPs on each LG ranged from 294 on LG 5 to 772 on LG 4. In sweet cherry, chromosome lengths ranged from 56.8 cM (LG 5) to 141.2 cM (LG 1), with a total map length of 655.4 cM. The number of SNPs on each LG ranged from 137 on LG 5 to 350 on LG 1.

3. Determining and resolving errors for haploblocks and haplotypes. The genetic maps of apple, peach, and sweet cherry were at first divided in 840, 103, 132 haploblocks, respectively, within which no recombination was observed in selected germplasm. After haplotype generation, 1262, 2012, and 74 Mendelian-inconsistent errors were reported by the `mconsistency.csv` file generated by FlexQTL. An additional 124, 429, and 64 recombinations were detected within the haploblocks for selected germplasm, resulting in the generation of additional haploblocks. The remaining Mendelian-inconsistent errors were mostly due to missing data within a haplotype that could not be resolved automatically. This missing data within haplotypes led to the assignment of haplotype numbers that were different to parental haplotypes that were therefore perceived as errors. In addition, some inconsistencies between SNP data and haplotype data were observed after haplotype generation that were easily resolved by looking at the ‘SNP Graph’ in GenomeStudio and adjusting either the haplotype or the SNP call. Finally, 33, zero, and zero SNPs were removed from the apple, peach, and sweet cherry data sets, respectively, because genotypes could not be accurately determined (Table 3).

The final number of haploblocks was 964, 135, and 196 for apple, peach, and sweet cherry respectively. For apple, the genetic length of the haploblocks varied between 0 and 7.77 cM with an average of 0.3 cM, the haploblocks contained between 1 and 15 SNPs, and the haploblocks contained an average of 4 SNPs. The number of haploblocks per apple LG ranged from 42 on LG 6 to 79 on LG 15, with an average of 57 haploblocks per LG. In peach, the length of the haploblocks varied between 0 cM and 30.47 cM with an average of 5.8 cM, the haploblocks contained between 1 and 210 SNPs, and the haploblocks contained an average of 30 SNPs. The number of haploblocks per peach LG ranged from 7 on LG 5 to 37 on LG 4, with an average of 17 haploblocks per LG. For sweet cherry, haploblocks had an average length of 2.6 cM, with a minimum of 0 cM and a maximum of 15.0 cM. The average number of SNPs per sweet cherry

haploblock was 8, with a minimum of 1 and a maximum of 61 SNPs. The average number of haploblocks per sweet cherry LG was 24, with a minimum of 16 haploblocks on LG 5 and LG 7 and a maximum of 47 haploblocks on LG 1.

SNP classification system. The final number of SNPs in the haplotyped data set was 3858, 4005, and 1617 for apple, peach, and sweet cherry, respectively. A total of 3350 (87%), 4005 (100%), and 1610 (99.6%) of these SNPs were classified as type 1 SNPs, which ultimately needed editing for less than 5% of their genotype calls in apple, peach, and sweet cherry, respectively (S8–S10 Tables). Type 2 SNPs, for which genotype clusters were shifted, totaled 300 (8%), zero, and seven (0.4%) SNPs for apple, peach, and sweet cherry, respectively, and this shift in cluster position lead to incorrect identification of one of the three clusters in the original automatic clustering by GenomeStudio. Type 3, SNPs with additional clusters, were assigned to 80 (2%), zero, and zero SNPs in apple, peach, and sweet cherry, respectively, and this presence of additional clusters led to incorrect genotype scoring of these SNPs that required subsequent curation. Type 4, SNPs with null alleles, were assigned to for 125 (3%), 145 (excluded from the final data set), and 43 (excluded from the final data set) SNPs in apple, peach, and sweet cherry, respectively, and these null alleles prevented correct automatic scoring for some individuals.

Discussion

We established a workflow to efficiently and confidently identify and remove genotyping errors from genotyped and pedigreed germplasm sets for apple, peach, and sweet cherry. The proposed workflow (Fig 1, S4 File) enables directed identification of markers and individuals with genotyping errors. It uses simple genetic principles such as inheritance of parental alleles, the co-segregation of linked markers, and the likelihood of double recombinations to find these errors. The order of steps was determined to efficiently minimize errors found in later steps and thereby minimize overall time needed to find errors in the data set. For example, in apple, any incorrect PC relationship would lead to an average of 196 reported Mendelian-inconsistent errors, and any unresolved Mendelian-inconsistent errors led to an average of 30 more reported Mendelian-consistent errors. The developed workflow was demonstrated on Illumina SNP array data and some software is specific to this platform (Table 4), but the same workflow order and genetic principles are appropriate for other marker types and genotyping platforms. The workflow is especially useful when medium- and high-throughput genotyping tools are used for which checking each individual marker would be too time-consuming.

Order and considerations of workflow steps

Different types of errors can be present in genotypic and pedigree data, caused by different kinds of issues (Table 1). To minimize the time needed for curation of these data, the proposed error checks need to be performed in a specific order. By first tackling issues that are common for many types of errors, subsequent curation of remaining errors becomes easier and quicker.

Removing individuals with low quality or irregular number of chromosome sets. The B-allele frequency plots provided a quick and easy way to identify and remove individuals with an irregular number of chromosome sets (polyploids and aneuploids) and individuals with low DNA quality. Removal of such individuals improved SNP calling and thus reduced the number of errors to be dealt with in later steps. A couple of individuals with poor quality that were originally kept, because of their importance as breeding parents, resulted in many PC errors. Making all their original SNP calls missing enabled automated imputation of most of these data points based on genetic information of relatives. Subsequent re-genotyping of these individuals matched the imputed data completely, confirming that the errors observed were

Table 4. Recommended software for each step of the genetic marker data curation workflow 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

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

<https://doi.org/10.1371/journal.pone.0210928.t004>

due to low-quality DNA samples and not to incorrect PC relationships. Polyploid and aneuploid individuals did not show a higher number of P(P)C errors, as expected. In contrast, these chromosome number abnormalities led to higher rates of false double recombination, either genome-wide (polyploid) or local [(segmental) aneuploids], that cannot be readily resolved other than by removal of these specific individuals.

The histogram function in GenomeStudio enabled quick identification of polyploids and individuals with very poor DNA samples without the need for additional steps in Excel, R, or other software. However, identification of aneuploids and individuals with potentially low-quality DNA samples was not as straightforward. Plotting the B-allele frequency against physical or genetic marker order (when available) required additional data manipulation and generation of the plots in software outside GenomeStudio, but most of it could be automated using R and custom scripts. Therefore, we suggest using GenomeStudio for initial removal of poor-quality samples and polyploids, and afterwards, when positional information for the markers is available, screening for aneuploids with the method described by Chagné and coworkers (2015) [52] (Table 4).

Obtaining a set of reliable SNPs

SNPs with major scoring issues that cannot be easily resolved manually need to be removed from the data set. The early detection and removal of these unreliable SNPs greatly reduces the number of marker and map errors reported, as well as the time spent evaluating these SNPs in later workflow stages. By using ASSIST, a quick subset of SNPs with robust genotype calls could be generated. On average across the three crops, 80% of this subset was retained in the final data set, which is lower than the 99% for single full-sib families that was reported by Di Guardo and co-workers (2015) [36]. As the number of generations and full-sib families in the germplasm increase, more SNPs with null alleles are likely to be detected and the more

complicated the genotype calling of these SNPs can become. In turn, this can lead to an increased discarding of SNPs, which could explain the lower proportion of SNPs retained in our germplasm sets compared to that reported by Di Guardo and co-workers (2015) [36].

Markers with null alleles identified by ASSIsT were removed from the data set, as they could only be identified and automatically called in specific F_1 families rather than in all families and across generations. However, many SNPs with null alleles that were later identified in the workflow could be accurately genotyped manually as long as homozygous 'AA' and 'BB' individuals could be distinguished from individuals that carried a null allele. This distinction was time-consuming and therefore we recommend saving these SNPs only when it justifies the time needed to do so. Examples when such markers can be of high value are in the construction of genetic linkage maps, even if multiple mapping populations are used [49], when they occur in a region of low coverage, or when they occur in a region of specific interest and help define additional alleles.

Very few other options exist to create a subset of high-quality genome-wide markers across pedigreed germplasm. GenomeStudio does provide several quality scores that have been used before in SNP filtering, but no guidelines exist on what threshold values to use. Using parameter thresholds regularly reported in literature [27,56–59] (GenTrain Score > 0.7, 50GC Score > 0.4, ClusterSep Score > 0.25, Call Rate > 0.9, and Minor Freq > 0.01) on the current data, the proportion of retained, unreliable, or monomorphic SNPs would be 12.3%, 23.1%, and 6.7% in apple, peach, and sweet cherry, respectively, and a large proportion of good SNPs would be discarded (27.8%, 28.2%, and 7.6%, respectively). Thus, ASSIsT greatly increased the number of reliable SNPs that were retained without reducing the quality of the subset of SNPs, making it the most efficient method to choose SNPs without prior knowledge on SNP performance (Table 4).

Updating pedigree records

As thresholds to confirm or discard historic pedigree information depends on the germplasm, genotyping platform, and data quality, they need to be assessed case-wise. A custom R-script provided quick and easy determination of the number of PC and PPC errors. However, the custom code required a significant amount of time to identify possible parents when one or both parents were unknown, especially for larger data sets. Similar issues were observed for Cervus, which took a long time to run (days) and did identify some incorrect relationships, especially for inbred material. Cervus also requires a specific data format and we experienced some problems running the software for large data sets that were not immediately resolved. GenomeStudio provided the quickest way to determine the number of PC and PPC errors, which could be determined immediately after loading the raw intensity data. However, new PC relationships could not automatically be determined and only SNPs retained by ASSIsT should be used when using GenomeStudio to determine the number of PC errors, to avoid inflating the number of PC errors. Therefore, we recommend using GenomeStudio to confirm existing pedigree records when using Illumina arrays and using an R-script to determine new, previously unknown, PC relationships (Table 4). Time-consuming analyses in R could be resolved by using a subset of markers equally spread across the genome. For confirming and identifying possible grandparent-grandchild relationships, we recommend the Excel template provided by van de Weg and co-workers (2018) [23]. However, this method can misconstrue aunts-uncles/nephew-nieces and individuals with other close relationships to the target individual as grandparents. Therefore, we recommend to only use this strategy when the user has a good understanding of the germplasm such as the origin of the material and the degree of inbreeding.

Individuals with only one parent known can still be used in a pedigree-based approach to find errors in the data set, although some errors might remain unnoticed. We recommend using the 'M_' and 'F_' prefixes to the individual's name to designate the unknown mother or father, respectively. When it is unclear whether the unknown individual is the mother or the father, the 'UP_' prefix can be used. Using this system instead of a non-descriptive name such as 'dummy 1' creates a clear connection between the individual with an unknown parent and the placeholder individual that is introduced. When the correct parent is later found, it also allows the quick replacement of the placeholder by the correct name (and corresponding genotypic data). Use of the same name for any missing parent should be avoided (e.g., using 'dummy' for all missing parents) unless the missing parent is unequivocally the parent of multiple individuals. If the same name is used incorrectly for multiple missing parents, the genotype of that single missing parent is expected by FlexQTL to be consistent with inheritance principles for all of its assigned offspring, potentially creating a large number of errors in further steps.

Although non-diploid individuals should be removed from the workflow before identifying reliable SNPs, they can have their pedigree checked if needed. Regardless of their ploidy, individuals should only contain alleles that are present in their parents. For example, a triploid individual with a marker call at one SNP of 'AAA' will be scored as 'AA', but can still not have a 'BB' parent. However, caution is advised as the grandparents through the parent that provided the unreduced gamete will also share a full allele set with any polyploid individual and thus these grandparents could also be incorrectly assigned as a parent of the polyploid individual. For example, the triploid 'Zonga' and its (diploid) grandparent 'Cox's Orange Pippin' share a full allele set (through an unreduced gamete of 'Alkmene') and thus no PC errors are reported [60]. However, only the combination of 'Delcorf' and 'Alkmene' could explain the genotypes of the triploid 'Zhonga' (AB+AA-AA test [23]). Thus, for triploids, not only do parents and offspring lead to no PC errors but some grandparents do as well, and the second parent is needed to identify the true PC relationship.

Creating or extending genetic maps

This study used available genetic maps for apple and cherry (i.e., [20,21,49,51]), integrated them when needed, and used available physical information (from [50] and [48]) to add any markers that were not already mapped. Some of these added markers were positioned at chromosome ends, which resulted in the increase of the map size by 7 cM for apple. In addition, the orientation of apple chromosome 5 was inverted here to match the orientation of the latest genome version [50]. If no genetic map is available, one will need to be constructed alongside genotypic data curation. The need for a precise genetic position of markers on the 9K peach array prompted development of consensus linkage map for peach [61] that in the future could serve as a reference map to estimate genetic positions of unmapped markers. A mapping approach for pedigreed, multi-parental maps is described by Di Pierro and co-workers (2016) [49].

Resolving remaining Mendelian-inconsistent errors

Use of GenomeStudio for detecting Mendelian-inconsistent errors is limited to Illumina array SNPs and cannot be used for other markers or haplotypes created in later steps of the workflow. In addition, some SNPs had their SNP scoring improved with ASSiST and manual curation, and thus the genotype scoring of GenomeStudio might not reflect the actual data. Although this latter limitation is also true when confirming pedigree data, the few differences in genotype calls between GenomeStudio and ASSiST are not expected to alter the outcome of

pedigree confirmation. In contrast, when resolving single Mendelian-inconsistent errors, it is important to know that the error is indeed present in the data set. Although Cervus counts the number of Mendelian-inconsistent errors, it does not report which markers are causing issues for which individuals, making it impractical to use to remove the remaining PC and PPC errors. In contrast to GenomeStudio, FlexQTL can handle multiple allele formats and is thus suited for the curation of both SNP data and haplotype data. In addition, FlexQTL checks for consistency over multiple generations, which enables detection of errors even if a genotype is missing in an intermediate individual. It also imputes missing data whenever possible. A disadvantage of FlexQTL is that it only reports one of the two individuals, often the parent, for which an error occurred; it is then up to the user to find the second individual, often the offspring, involved in the Mendelian-inconsistent error. Therefore, we recommend using FlexQTL to identify Mendelian-inconsistent errors and resolving them with the help of GenomeStudio (Table 4).

Using map and phasing information to detect Mendelian-consistent errors

FlexQTL performed very accurate phasing and only a few phasing issues were noticed. Most of these phasing issues were observed as double recombinations in offspring of an individual that served as a founder. The lack of parental info for this founder provided FlexQTL more freedom to phase alleles, as the phasing in the founder did not need to match its parents. Incorrect phasing was most likely caused by one or very few offspring for which a true recombination occurred in the map region. In those individuals, no double recombination occurred, and the incorrect phasing inferred by FlexQTL minimized the interval over which the true recombination occurred. However, this minimalization of the recombination interval incorrectly specified where the recombination had occurred, causing incorrect phasing and resulting in one or multiple false double recombinations in full- and half-sibs of the individual(s) with the true recombination. Making genotype calls missing for the individual(s) with a recombination in that area enlarged the recombination interval for those individuals, but also led to correct phasing in their parent and resolved the supposed double recombinations in their full- and half-sibs. Very few other phasing issues were observed that could not be resolved on a single SNP level but were later resolved at the haploblock level. Thus, a small number of phasing issues can be accepted when moving forward to generating haploblocks and they could be nullified by FlexQTL by setting the parameter 'DeleteDR' to 1.

Haploblock and haplotype determination

Visual FlexQTL showed good accuracy (between 12% and 33% of the initial haploblocks had to be divided into additional haploblocks to avoid recombination within haploblocks for selected material) in determining haploblock borders based on historic recombination events. Two reasons exist for not identifying all historic recombinations for haploblock border determination. First, Visual FlexQTL determines the border as the middle of the recombination interval. The more non-informative markers present in the recombination interval (due to homozygosity or lack of co-segregation (phase) information), the less likely that the middle position is the true position of the historic recombination (which determines the haploblock border). Secondly, FlexQTL determines haploblock borders sequentially, starting with small recombination intervals; if multiple recombinations occur in the same region, one haploblock border could suffice to account for all recombinations. This approach thus minimizes the number of recombination sites needed to explain observed segregation data. In reality, the recombinations could have occurred between different markers, requiring that region to be

split in additional haploblocks to avoid recombination within haploblocks for selected material.

PediHaplotyper's haplotypes did not always match with SNP data. In most cases, these inconsistencies were introduced during the marker consistency check with FlexQTL to ensure the haplotypes in an individual matched those of its parents and offspring. When the haplotype that caused the inconsistency was represented well in the pedigree, the haplotype was correct and the original genotype call for the SNP was incorrect. Thus, in these cases, haplotype curation identified additional errors in the SNP data. These errors were mostly caused by (very) small incorrectly identified genotype clusters or by single calling errors in the data set that were not detected earlier. When haplotypes in poorly represented individuals (one or two directly related individuals in the data set) showed an inconsistency with the SNP data, the SNP data was mostly correct and an error had occurred during haplotyping. The error could span multiple generations leading to inconsistencies for multiple individuals but its impact on the data set was small as the overall representation of the incorrect haplotype was small. In the rare case that a poor representation led to incorrect haplotype determination, the actual cause of the inconsistency often remained unclear, but for some it was due to a recombination within a haploblock for an un-genotyped ancestor or one of the direct parents of such an ancestor.

Haploblock borders are not fixed and can change based on the application of the final data set and the germplasm used. For example, for QTL analyses some of the haploblocks defined here will be too large as they span multiple cM; they will show within-haploblock recombination in numerous unselected offspring thereby increasing the number of missing haplotype calls thus increasing uncertainty in QTL position (including the widening of QTL intervals). Haploblock sizes can therefore be reduced to minimize within haploblock recombination and better define QTL regions. However, when haploblocks are very small, many haploblocks will consist of only one SNP or a few SNPs, increasing data sizes (and thereby computation time in downstream analyses) and reducing the number of haplotypes per haploblock, which can reduce the suitability of the data for visual examination. Unlike the 8K apple SNP array, the 20K apple SNP array was designed to have clusters of multiple SNPs spread at approximately 1 cM intervals. A similar approach was used to create 9K add-ons for the 9K peach array and 6K cherry array [62]. This strategy supports the generation of haploblocks consisting of SNPs aggregated within 1 cM intervals while still having multiple SNPs in a single haploblock and thus multiple informative haplotypes.

Different germplasm will also lead to different haploblock borders. Currently, haploblocks are based on historic recombination events representing the U.S. breeding programs included in this study. Other breeding programs or genetic studies might have other sets of founders and thus different recombinations of relevance. Furthermore, the addition of new advanced selections and parents will introduce new recombinations in their germplasm. Finally, as the understanding of the apple, peach, and cherry germplasm increases, previously unknown progenitors, founders, and pedigree connections will be discovered, also increasing the number of observed recombinations.

Given that haploblocking is performed at a relative late stage in the workflow, haploblock borders can be altered without the need to redo all previously conducted pedigree and SNP marker curation. In fact, existing haplotype data can be converted back to phased, fully curated SNP data which, in turn, can be used to determine haplotypes for any set of haploblocks. As the SNPs are already phased and missing SNP data was imputed based on the haplotypes, haplotype determination for new haploblock borders should not create new genotyping errors in the data set. Once numbers of new recombinations are high enough to justify updating of haploblock data, part of the haploblocks and their haplotypes should be altered. PediHaplotyper supports the use of previous haplotype definitions for haploblocks that did not change in

composition. Adjusted haploblocks could be marked through their names, thus providing tools to monitor new as well as previous, possibly well-known, marker alleles.

The SNP classification system and integration of genotypic data for new germplasm into existing data sets

The established SNP classification system enables the quick creation of a subset of SNPs that require minimal or no data curation and provides a guideline on possible issues with other SNPs and how to solve them. The system should help with the quick integration of new genotypic data into existing data sets. Genotype calls for SNPs of type 1 and type 2 can be quickly integrated with high confidence in their genotype calls. Where desired, SNPs of type 3, 4, and 5 can also be integrated, but additional curation would be required. Depending on germplasm tested, these SNPs might have incorrect genotype scoring but their SNP type is an indication of why the genotype scoring is wrong and how to fix it. In other germplasm, additional SNPs in the probe or null alleles might not be present, causing SNPs that are now classified as type 4 or type 5 to give reliable results as if they were type 1 or type 2. Similarly, if germplasm is used that is unrelated to that used here, type 1 and type 2 SNPs might show additional clusters or null alleles and will require further curation. Finally, type 7 SNPs, which could not be mapped in this germplasm, might be mapped and valuable for other germplasm.

The available reference data (www.rosaceae.org), combined with the SNP classification system, will facilitate correct curation of additional genotypic data, even if the new germplasm is not directly descended. The SNP genotype calls provided here are a reference for the genotype of each observed genotype cluster in GenomeStudio. In addition, SNP cluster coordinates of the latest GenomeStudio file can be imported into new projects, thus helping GenomeStudio to correctly identify clusters. Finally, the use of reference iScan data is especially useful for markers that have only two of the three clusters in a new project but all three clusters were defined in the current reference data set. By adding reference iScan data into the new project, all three clusters will be available, ensuring correct automated genotype calling. Therefore, we recommend including available reference data when obtaining genotype calls for new germplasm.

Data curation in apple

The need for SNP data curation in apple was increased by the whole genome duplication in the evolutionary history of apple, the relatively poorer quality of the first genome draft used for development of the 8K SNP array, and unidentified polymorphisms in the probe regions during SNP array design. The genome duplication resulted in presence of multiple highly similar sequences on different chromosomes. Indeed, a BLAST analysis against *Malus* genome v1.0 of the first 24 nucleotides of the 3' region of arrayed SNP probes, which is most important for probe binding, showed that approximately 50% of the sequences returned multiple hits with almost all of these hits being located on multiple LGs [34]. This proportion is expected to be lower for the latest genome version [50] as most errors in assembly were removed but the proportion is expected to remain high due to chromosome and gene duplication observed in apple. Where two genomic regions are targeted by the same probe, complex cluster plots will occur if more than one of the targeted loci segregate within a single family. Such markers must be excluded from a curated data set. Multi-target markers might still be robust if they segregate at only one locus. In this case, only the cluster plot space is reduced (mostly halved), causing clusters to be located more closely to each other. In turn, this might occasionally cause separation issues. Also, some markers are lost because GenomeStudio cannot assign genotype calls for markers where one of the homozygous clusters is located at $\theta = 0.5$, the center of the x-

axis, and thus these markers are considered by the software to have failed. A special case for two-locus markers occurs where each locus segregates in specific families but both loci never segregate together in the same family. In this case, genotype scoring might be performed accurately, and the SNP still needs to be present twice in the map although under different names. Two- and three-locus SNPs have been successfully mapped in the multi-family based genetic linkage map created by Di Pierro and co-workers (2016) [49]. However, in subsequent QTL mapping studies on pedigreed germplasm, such markers were excluded, as in the current study.

Several intermediate progenitors in the apple data set lacked any genotypic data and therefore the recorded link between some important breeding parents and their ancestors had to be set to unknown during haploblock and haplotype determination. For some progenitors, 20K data from the European FruitBreedomics project was available that reestablished the connection between genotyped individuals and their ancestors, but many other progenitors likely no longer exist. Individuals that were disconnected from the pedigree with little representation could not therefore have their haplotypes accurately determined using PediHaplotyper. It was, however, possible to manually determine their haplotypes based on their SNP data and haplotypes present in disconnected relatives.

Data curation in peach

In peach, the most challenging step in the workflow was the curation of pedigree information over nine generations. Although much pedigree information is available in the literature [63], we identified incorrect parentage in the PC error analysis in cultivars and breeding selections, which we attributed to selfing or outcrossing. Incorrect pedigree records were previously reported in the UC Davis processing peach breeding program in approximately 20% of individuals, both parental and breeding selections [16]. In this work, we identified incorrect parentage in approximately 11% of the pedigree records from the three fresh market peach breeding programs, most of which were observed in breeding selections. High level of inbreeding and coancestry in the U.S cultivated peach germplasm [64] creates overlap in the ancestral generations of most U.S. peach breeding programs. Therefore, corrections in the ancestral pedigree records reported by Fresno-Ramírez and co-workers (2015) [16] reduced the number of errors detected here. Furthermore, intermediate parents were unavailable for genotyping, so pedigree connections were preserved by retaining pedigree information even though many intermediate progenitors were not genotyped. Finally, the presence of missing data within a haplotype resulted in Mendelian-inconsistent errors in the haploblock and haplotype generation steps, which made the haploblock data curation time-consuming.

Data curation in sweet cherry

For the sweet cherry germplasm, the most challenging issue was the small sample size of some families (as few as four individuals), which were too small for FlexQTL to accurately determine linkage phase. For parents with just one genotyped offspring, phasing of the parent homologs was considered putative as recombination inherited by offspring could not be determined. For those parents with just two genotyped offspring, recombinations were arbitrarily assigned between the two offspring, as the true recombinant offspring could not be determined. In addition, scarce information on pedigrees in ancestral generations beyond about five limited further imputations in data curation, unlike for apple and peach. Various founders showed extensive regions of common haplotypes, indicating a high degree of relatedness among such founders. Some recently published haplotyping results exemplify this for the founders 'Black Republican' and 'Napoleon' [21]. Unraveling the unknown relationships among founders

could thus provide useful information for future data curation in sweet cherry. Finally, the use of the peach genome to determine genetic position for the sweet cherry SNPs did not seem to result in larger map issues compared to the other crops: no SNPs were moved to different positions in the map during data curation and a similar proportion of SNPs were discarded due to map issues for sweet cherry as for peach (~4%). However, the use of the peach genome might have negatively affected development of the cherry SNP array as a much lower proportion of the SNPs on the array were retained in the final sweet cherry data set compared to the other crops.

Expectations for other crops

The proposed workflow could be applied to other diploid crops with similar breeding systems where clonally propagated relatives of current breeding material still exist. However, there are additional aspects that would need to be considered in certain circumstances that were not encountered in the present study. First, this workflow makes the assumption that there are no differences in the true SNP map order among individuals of a species. In interspecific crosses where there can be differences in chromosome arrangements between parental species, the different SNP order or indel variation among individuals could result in additional perceived double recombinations or other difficulties in following this workflow. Additionally, this workflow assumes that there is sufficient marker information to correctly identify pedigree relationships and assumes sufficient segregation information for validating marker order and identifying Mendelian-consistent errors. When using highly homozygous, inbred individuals, there might be too few segregating markers available to correctly identify marker order or find Mendelian-consistent errors through double recombinations. Also, for small germplasm sets, too few recombinations might be available to detect incorrect marker order. Finally, the prevalence of missing genotypic values should be sufficiently low across individuals. Unlike the SNP arrays used in this study, some genotyping methods such as Genotyping-by-Sequencing do not consistently target specific loci. This non-specificity can increase the flexibility of their use, but also raises new issues for which the current workflow would have to be adapted, including the potential decrease in accuracies of genotyping and haploblock determination due to unbalanced representation of genotyped loci, high levels of missing data, and sequencing errors.

High-quality archived SNP and haplotype data sets

The presented genome-wide genotypic data sets for apple, peach, and sweet cherry are of very high quality, are composed of genetically complex germplasm, and contain no errors that could be determined based on pedigree information. This high quality provides confidence in the results of downstream analyses. Such confidence is important as many of these results are expected to lead to fundamental discoveries and practical breeding application. The iScan data, phased SNP, and haplotype data sets of individuals in the apple, peach, and sweet cherry crop reference sets are available through the Genome Database for Rosaceae (www.rosaceae.org).

Marker and pedigree data from germplasm subsets of the current U.S. RosBREED project, the EU-FruitBreedomics project, and other research projects have previously been curated by a precursor to the current workflow and used for the creation of a multi-family based genetic linkage maps [20,49] and in multifamily based QTL studies in apple [65–68], peach [22,69], and sweet cherry [14]. Also, elements of the workflow were used for allo-octoploid strawberry to curate Axiom-based SNP markers [32] and pedigree data that were subsequently used in multi-family based QTL analyses [70–72]. While providing high-quality data for each analysis separately, these earlier steps in data curation have helped guide and streamline the data

curation workflow presented here. The current workflow and resulting data sets ensure that the same curation steps have been used across the data sets of multiple crops and that the data sets are of the same high quality.

Conclusion

A curation workflow for genotypic data of pedigreed germplasm was generated by determining the optimal order of resolving issues and by providing a step-by-step guideline. Using simple genetic principles, errors can be found and curated in a directed and efficient way, reducing the time needed to obtain a high-quality genotypic data set. The workflow was used to obtain a SNP data set for large germplasm sets for each of apple, peach, and sweet cherry representing U.S. breeding programs based on the apple 8K SNP array, peach 9K SNP array, and cherry 6K SNP array, respectively, whose SNP data is available through this paper (www.rosaceae.org), as well as used on apple and peach germplasm sets representing European breeding programs based on the apple 20K and peach 9K arrays, whose SNP data are still private. These high-quality data sets contain the largest sets of SNPs obtained through their respective SNP arrays and will provide the foundation for confident subsequent analyses in genetic research.

Supporting information

S1 Table. Apple germplasm genotyped and used for data curation workflow. Individuals are split over the publicly available RosBREED Crop Reference Set, three privately held RosBREED Breeding Pedigree Sets, and genotypic data received from either KULeuven (Belgium) or the FruitBreedomics project. Except for the Breeding Pedigree sets, curated pedigree information is given for each individual. For each individual, the type of material (selected vs. unselected), the location of sampling, quality of the results, and inferred ploidy of the sample are given. For unselected seedlings, the family to which they belong is also given. For the Breeding Pedigree Sets, this information is summarized per full-sib family. If tissue was collected at the USDA germplasm repository in Geneva, a GRIN accession number is also provided. Parents highlighted in yellow did not have genotypic data and their pedigree-relationships could not be tested.
(XLSX)

S2 Table. Peach germplasm genotyped and used for curation workflow. Individuals are split over the publicly available RosBREED Crop Reference Set and three privately held RosBREED Breeding Pedigree Sets. Except for the Breeding Pedigree Sets, curated pedigree information is given for each individual. For each individual, the type of material (selected vs. unselected), the location of sampling, and quality of the results of the sample are given. For unselected seedlings, the family to which they belong is also given. For the Breeding Pedigree Sets, this information is summarized per full-sib family.
(XLSX)

S3 Table. Sweet cherry germplasm genotyped and used for curation workflow. All individuals are part of the publicly available RosBREED Crop Reference Set. For each individual, curated pedigree information, the type of material (selected vs. unselected), the location of sampling, and quality of the results of the sample are given. For unselected seedlings, the family to which they belong is also given.
(XLSX)

S4 Table. Parameter settings used for (A) filtering SNPs used in analyses of B-allele frequency, (B) running ASSiST, (C) running FlexQTL for detecting Mendelian-inconsistent errors and

Mendelian-consistent errors, and (D) running FlexQTL for phasing, haploblock determination, and creating PediHaplotyper input files.

(XLSX)

S5 Table. Final genetic map used for apple during data curation. For each marker, genetic position, associated haploblock, and physical position based on the apple GDDH 13 v1.1 genome are given.

(XLSX)

S6 Table. Final genetic map used for peach during data curation. For each marker, genetic position, associated haploblock, and physical position based on the peach v2 genome are given.

(XLSX)

S7 Table. Final genetic map used for sweet cherry during data curation. For each marker, genetic position, associated haploblock, and physical position based on the peach v2 genome are given.

(XLSX)

S8 Table. SNP classification for apple. Each SNP is classified as follows: Type '1' for SNPs with good clustering and less than 5% call errors, '2' for SNPs with shifted clusters causing one of the clusters to be called incorrectly, '3' for SNPs with additional clusters (excluding null-alleles) that cause the incorrect identification of at least one cluster, '4' for SNPs with null-alleles that cannot be correctly called automatically, '5' for SNPs that could not be mapped accurately but had correct clustering, '6' for monomorphic SNPs, and '7' for failed SNPs.

(XLSX)

S9 Table. SNP classification for peach. Each SNP is classified as follows: Type '1' for SNPs with good clustering and less than 5% call errors, '2' for SNPs with shifted clusters causing one of the clusters to be called incorrectly, '3' for SNPs with additional clusters (excluding null-alleles) that cause the incorrect identification of at least one cluster, '4' for SNPs with null-alleles that cannot be correctly called automatically, '5' for SNPs that could not be mapped accurately but had correct clustering, '6' for monomorphic SNPs, and '7' for failed SNPs.

(XLSX)

S10 Table. SNP classification for sweet cherry. Each SNP is classified as follows: Type '1' for SNPs with good clustering and less than 5% call errors, '2' for SNPs with shifted clusters causing one of the clusters to be called incorrectly, '3' for SNPs with additional clusters (excluding null-alleles) that cause the incorrect identification of at least one cluster, '4' for SNPs with null-alleles that cannot be correctly called automatically, '5' for SNPs that could not be mapped accurately but had correct clustering, '6' for monomorphic SNPs, and '7' for failed SNPs.

(XLSX)

S1 Fig. Call rates observed for individuals classified as having good, intermediate, or bad quality of genotypic data as defined by their B-allele frequency plot outcome. Higher call rates are observed for individuals with better quality of genotypic data.

(DOCX)

S2 Fig. SNP B-allele frequencies plotted against physical position in the genome for (A) triploid individuals excluding 'Jonagold', and (B) individuals with a tetraploid pattern.

(DOCX)

S1 File. Description of procedures performed during the steps of the workflow.
(DOCX)

S2 File. R-script used to create B-allele frequency plots for all genotyped individuals.
(R)

S3 File. R-scripts used to confirm and deduce P(P)C relationships.
(R)

S4 File. Hands-on guideline on how to perform data curation using the steps described in this study.
(PDF)

Acknowledgments

The authors thank Jasper Rees (Agricultural Research Council, South Africa) for help in genotyping apple samples.

Author Contributions

Conceptualization: Eric Van de Weg, Cameron Peace.

Data curation: Stijn Vanderzande, Nicholas P. Howard, Lichun Cai, Cassia Da Silva Linge, Laima Antanaviciute, Cameron Peace.

Formal analysis: Stijn Vanderzande, Nicholas P. Howard, Lichun Cai, Cassia Da Silva Linge, Laima Antanaviciute.

Funding acquisition: Nahla Bassil, Ksenija Gasic, Amy Iezzoni, Eric Van de Weg, Cameron Peace.

Methodology: Stijn Vanderzande, Nicholas P. Howard, Lichun Cai, Cassia Da Silva Linge, Marco C. A. M. Bink, Johannes W. Kruisselbrink, Eric Van de Weg, Cameron Peace.

Project administration: Amy Iezzoni, Cameron Peace.

Resources: Nahla Bassil, Ksenija Gasic, Amy Iezzoni, Eric Van de Weg, Cameron Peace.

Software: Stijn Vanderzande, Nicholas P. Howard, Marco C. A. M. Bink, Johannes W. Kruisselbrink, Eric Van de Weg.

Supervision: Ksenija Gasic, Amy Iezzoni, Eric Van de Weg, Cameron Peace.

Writing – original draft: Stijn Vanderzande, Nicholas P. Howard.

Writing – review & editing: Lichun Cai, Cassia Da Silva Linge, Marco C. A. M. Bink, Johannes W. Kruisselbrink, Ksenija Gasic, Amy Iezzoni, Eric Van de Weg, Cameron Peace.

References

1. Pompanon F, Bonin A, Bellemain E, Taberlet P. Genotyping errors: causes, consequences and solutions. *Nat Rev Genet.* 2005; 6: 847–859. <https://doi.org/10.1038/nrg1707> PMID: 16304600
2. Buetow KH. Influence of aberrant observations on high-resolution linkage analysis outcomes. *Am J Hum Genet.* 1991; 49: 985–994. PMID: 1928104
3. Goldstein DR, Zhao H, Speed TP. The effects of genotyping errors and interference on estimation of genetic distance. *Hum Hered.* 1997; 47: 86–100. <https://doi.org/10.1159/000154396> PMID: 9097090
4. Hackett CA, Broadfoot LB. Effects of genotyping errors, missing values and segregation distortion in molecular marker data on the construction of linkage maps. *Heredity.* 2003; 90: 33–38. <https://doi.org/10.1038/sj.hdy.6800173> PMID: 12522423

5. Abecasis GR, Cherny SS, Cardon LR. The impact of genotyping error on family-based analysis of quantitative traits. *European Journal of Human Genetics*. 2001; 9: 130–134. <https://doi.org/10.1038/sj.ejhg.5200594> PMID: 11313746
6. Gordon D, Finch SJ, Nothnagel M, Ott J. Power and sample size calculations for case-control genetic association tests when errors are present: application to single nucleotide polymorphisms. *Hum Hered*. 2002; 54: 22–33. <https://doi.org/10.1159/000066696> PMID: 12446984
7. Browning BL, Yu Z. Simultaneous genotype calling and haplotype phasing improves genotype accuracy and reduces false-positive associations for genome-wide association studies. *Am J Hum Genet*. 2009; 85: 847–861. <https://doi.org/10.1016/j.ajhg.2009.11.004> PMID: 19931040
8. Terwilliger J, Weeks D, Ott J. Laboratory errors in the reading of marker alleles cause massive reductions in LOD score and lead to gross overestimates of the recombination fraction. *Am J Hum Genet*. 1990; 47: A201.
9. Kirk KM, Cardon LR. The impact of genotyping error on haplotype reconstruction and frequency estimation. *Eur J Hum Genet*. 2002; 10: 616–622. <https://doi.org/10.1038/sj.ejhg.5200855> PMID: 12357332
10. Cheung CYK, Thompson EA, Wijsman EM. Detection of Mendelian consistent genotyping errors in pedigrees. *Genet Epidemiol*. 2014; 38: 291–299. <https://doi.org/10.1002/gepi.21806> PMID: 24718985
11. Vouillamoz JF, Grando MS. Genealogy of wine grape cultivars: “Pinot” is related to “Syrah.” *Heredity (Edinb)*. 2006; 97: 102–110. <https://doi.org/10.1038/sj.hdy.6800842> PMID: 16721391
12. Evans KM, Patocchi A, Rezzonico F, Mathis F, Durel CE, Fernández-Fernández F, et al. Genotyping of pedigreed apple breeding material with a genome-covering set of SSRs: trueness-to-type of cultivars and their parentages. *Mol Breeding*. 2011; 28: 535–547. <https://doi.org/10.1007/s11032-010-9502-5>
13. Lacombe T, Boursiquot J-M, Laucou V, Di Vecchi-Staraz M, Péros J-P, This P. Large-scale parentage analysis in an extended set of grapevine cultivars (*Vitis vinifera* L.). *Theor Appl Genet*. 2013; 126: 401–414. <https://doi.org/10.1007/s00122-012-1988-2> PMID: 23015217
14. Rosyara UR, Sebolt AM, Peace C, Iezzoni AF. Identification of the paternal parent of ‘Bing’ sweet cherry and confirmation of descendants using single nucleotide polymorphism markers. *J Amer Soc Hort Sci*. 2014; 139: 148–156.
15. Pikunova A, Madduri M, Sedov E, Noordijk Y, Peil A, Troglio M, et al. ‘Schmidt’s Antonovka’ is identical to ‘Common Antonovka’, an apple cultivar widely used in Russia in breeding for biotic and abiotic stresses. *Tree Genetics & Genomes*. 2014; 10: 261–271. <https://doi.org/10.1007/s11295-013-0679-8>
16. Fresnedo-Ramírez J, Crisosto CH, Gradziel TM, Famula TR. Pedigree correction and estimation of breeding values for peach genetic improvement. *Acta Horticulturae*. 2015; 1084: 249–256. <https://doi.org/10.17660/ActaHortic.2015.1084.35>
17. Fresnedo-Ramírez J, Frett TJ, Sandefur PJ, Salgado-Rojas A, Clark JR, Gasic K, et al. QTL mapping and breeding value estimation through pedigree-based analysis of fruit size and weight in four diverse peach breeding programs. *Tree Genetics & Genomes*. 2016; 12: 25. <https://doi.org/10.1007/s11295-016-0985-z>
18. Lassois L, Denancé C, Ravon E, Guyader A, Guisnel R, Hibrand-Saint-Oyant L, et al. Genetic diversity, population structure, parentage analysis, and construction of core collections in the French apple germplasm based on SSR markers. *Plant Mol Biol Rep*. 2016; 34: 827–844. <https://doi.org/10.1007/s11105-015-0966-7>
19. Larsen B, Toldam-Andersen TB, Pedersen C, Ørgaard M. Unravelling genetic diversity and cultivar parentage in the Danish apple gene bank collection. *Tree Genetics & Genomes*. 2017; 13: 14. <https://doi.org/10.1007/s11295-016-1087-7>
20. Howard NP, van de Weg E, Bedford DS, Peace CP, Vanderzande S, Clark MD, et al. Elucidation of the ‘Honeycrisp’ pedigree through haplotype analysis with a multi-family integrated SNP linkage map and a large apple (*Malus domestica*) pedigree-connected SNP data set. *Horticulture Research*. 2017; 4: 17003. <https://doi.org/10.1038/hortres.2017.3> PMID: 28243452
21. Cai L, Voorrips RE, van de Weg E, Peace C, Iezzoni A. Genetic structure of a QTL hotspot on chromosome 2 in sweet cherry indicates positive selection for favorable haplotypes. *Mol Breeding*. 2017; 37: 85. <https://doi.org/10.1007/s11032-017-0689-6>
22. Hernández Mora JR, Micheletti D, Bink MAM, Van de Weg WE, Bassi D, Nazzicari N, et al. Discovering peach QTLs with multiple progeny analysis. *Acta Horticulturae*. 2017; 1172: 405–410. <https://doi.org/10.17660/ActaHortic.2017.1172.77>
23. van de Weg E, Di Guardo M, Jänsch M, Socquet-Juglard D, Costa F, Baumgartner I, et al. Epistatic fire blight resistance QTL alleles in the apple cultivar ‘Enterprise’ and selection X-6398 discovered and characterized through pedigree-informed analysis. *Mol Breeding*. 2018; 38: 5. <https://doi.org/10.1007/s11032-017-0755-0>

24. Hoffman JI, Amos W. Microsatellite genotyping errors: detection approaches, common sources and consequences for paternal exclusion. *Mol Ecol*. 2005; 14: 599–612. <https://doi.org/10.1111/j.1365-294X.2004.02419.x> PMID: 15660949
25. Xu X, Bai G. Whole-genome resequencing: changing the paradigms of SNP detection, molecular mapping and gene discovery. *Mol Breeding*. 2015; 35: 33. <https://doi.org/10.1007/s11032-015-0240-6>
26. Rasheed A, Hao Y, Xia X, Khan A, Xu Y, Varshney RK, et al. Crop breeding chips and genotyping platforms: progress, challenges, and perspectives. *Mol Plant*. 2017; 10: 1047–1064. <https://doi.org/10.1016/j.molp.2017.06.008> PMID: 28669791
27. Chagné D, Crowhurst RN, Troggio M, Davey MW, Gilmore B, Lawley C, et al. Genome-wide SNP detection, validation, and development of an 8K SNP array for apple. *PLOS ONE*. 2012; 7: e31745. <https://doi.org/10.1371/journal.pone.0031745> PMID: 22363718
28. Verde I, Bassil N, Scalabrin S, Gilmore B, Lawley CT, Gasic K, et al. Development and evaluation of a 9K SNP array for peach by internationally coordinated SNP detection and validation in breeding germplasm. *PLOS ONE*. 2012; 7: e35668. <https://doi.org/10.1371/journal.pone.0035668> PMID: 22536421
29. Peace C, Bassil N, Main D, Ficklin S, Rosyara UR, Stegmeir T, et al. Development and evaluation of a genome-wide 6K SNP array for diploid sweet cherry and tetraploid sour cherry. *PLOS ONE*. 2012; 7: e48305. <https://doi.org/10.1371/journal.pone.0048305> PMID: 23284615
30. Le Paslier M-C, Choisne N, Bacilieri R, Bounon R, Boursiquot J-M, Bras M, et al. The GrapeReSeq 18K Vitis genotyping chip. 9th International symposium grapevine physiology and biotechnology. La Serena, Chile; 2013. p. 123.
31. Bianco L, Cestaro A, Sargent DJ, Banchi E, Derdak S, Guardo MD, et al. Development and validation of a 20K single nucleotide polymorphism (SNP) whole genome genotyping array for apple (*Malus × domestica* Borkh). *PLOS ONE*. 2014; 9: e110377. <https://doi.org/10.1371/journal.pone.0110377> PMID: 25303088
32. Bassil NV, Davis TM, Zhang H, Ficklin S, Mittmann M, Webster T, et al. Development and preliminary evaluation of a 90K Axiom® SNP array for the allo-octoploid cultivated strawberry *Fragaria × ananassa*. *BMC Genomics*. 2015; 16: 155. <https://doi.org/10.1186/s12864-015-1310-1> PMID: 25886969
33. Bianco L, Cestaro A, Linsmith G, Muranty H, Denancé C, Théron A, et al. Development and validation of the Axiom® Apple 480K SNP genotyping array. *Plant J*. 2016; 86: 62–74. <https://doi.org/10.1111/tpl.13145> PMID: 26919684
34. Troggio M, Šurbanovski N, Bianco L, Moretto M, Giongo L, Banchi E, et al. Evaluation of SNP data from the *Malus* Infinium array identifies challenges for genetic analysis of complex genomes of polyploid origin. *PLOS ONE*. 2013; 8: e67407. <https://doi.org/10.1371/journal.pone.0067407> PMID: 23826289
35. Di Guardo M, Micheletti D, Bianco L, Koehorst-van Putten HJJ, Longhi S, Costa F, et al. ASSIsT: an automatic SNP scoring tool for in- and outbreeding species—Reference Manual. 2015.
36. Di Guardo M, Micheletti D, Bianco L, Koehorst-van Putten HJJ, Longhi S, Costa F, et al. ASSIsT: an automatic SNP scoring tool for in- and outbreeding species. *Bioinformatics*. 2015; 31: 3873–3874. <https://doi.org/10.1093/bioinformatics/btv446> PMID: 26249809
37. Voorrips RE, Bink MCAM, Krusselbrink JW, Koehorst-van Putten J, van de Weg WE. PediHaplotyper: software for consistent assignment of marker haplotypes in pedigrees. *Mol Breed*. 2016; 36. <https://doi.org/10.1007/s11032-016-0539-y> PMID: 27547106
38. Iezzoni A, Weebadde C, Luby J, Chengyan Yue, van de Weg E, Fazio G, et al. Rosbreed: enabling marker-assisted breeding in Rosaceae. *Acta Horticulturae*. 2010; 859: 389–394. <https://doi.org/10.17660/ActaHortic.2010.859.47>
39. Iezzoni A, Weebadde C, Peace C, Main D, Bassil NV, Coe M, et al. Where are we now as we merge genomics into plant breeding and what are our limitations? Experiences from RosBREED. *Acta Horticulturae*. 2016; 1117: 1–5. <https://doi.org/10.17660/ActaHortic.2016.1117.1>
40. Iezzoni A, Peace C, Main D, Bassil N, Coe M, Finn C, et al. RosBREED2: progress and future plans to enable DNA-informed breeding in the Rosaceae. *Acta Horticulturae*. 2017; 1172: 115–118. <https://doi.org/10.17660/ActaHortic.2017.1172.20>
41. Laurens F, Durel C-E, Patocchi A, Peil A, Salvi S, Tartarini S, et al. Review on apple genetics and breeding programs and presentation of a new initiative of a new European initiative to increase fruit breeding efficiency. *Journal of Fruit Science*. 2010; 27: 102–107.
42. Laurens F, Aranzana MJ, Arús P, Bassi D, Bonany J, Corelli L, et al. Review of fruit genetics and breeding programmes and a new European initiative to increase fruit breeding efficiency. *Acta Horticulturae*. 2012; 929: 95–102. <https://doi.org/10.17660/ActaHortic.2012.929.12>
43. Laurens F, Aranzana MJ, Arús P, Bassi D, Bink M, Bonany J, et al. An integrated approach for increasing breeding efficiency in apple and peach in Europe. *Hortic Res*. 2018; 5: 11. <https://doi.org/10.1038/s41438-018-0016-3> PMID: 29507735

44. Peace CP, Luby JJ, van de Weg WE, Bink MCAM, Iezzoni AF. A strategy for developing representative germplasm sets for systematic QTL validation, demonstrated for apple, peach, and sweet cherry. *Tree Genetics & Genomes*. 2014; 10: 1679–1694. <https://doi.org/10.1007/s11295-014-0788-z>
45. Velasco R, Zharkikh A, Affourtit J, Dhingra A, Cestaro A, Kalyanaraman A, et al. The genome of the domesticated apple (*Malus × domestica* Borkh.). *Nat Genet*. 2010; 42: 833–839. <https://doi.org/10.1038/ng.654> PMID: 20802477
46. Peace C, Bianco L, Troggio M, Van de Weg E, Howard NP, Cornille A, et al. Apple whole genome sequences: recent advances and new prospects. *Horticulture Research*. 2019; *Accepted*.
47. Verde I, Abbott AG, Scalabrin S, Jung S, Shu S, Marroni F, et al. The high-quality draft genome of peach (*Prunus persica*) identifies unique patterns of genetic diversity, domestication and genome evolution. *Nat Genet*. 2013; 45: 487–494. <https://doi.org/10.1038/ng.2586> PMID: 23525075
48. Verde I, Jenkins J, Dondini L, Micali S, Pagliarani G, Vendramin E, et al. The Peach v2.0 release: high-resolution linkage mapping and deep resequencing improve chromosome-scale assembly and contiguity. *BMC Genomics*. 2017; 18: 225. <https://doi.org/10.1186/s12864-017-3606-9> PMID: 28284188
49. Di Pierro EA, Gianfranceschi L, Di Guardo M, Koehorst-van Putten HJ, Kruisselbrink JW, Longhi S, et al. A high-density, multi-parental SNP genetic map on apple validates a new mapping approach for outcrossing species. *Horticulture Research*. 2016; 3: 16057. <https://doi.org/10.1038/hortres.2016.57> PMID: 27917289
50. Daccord N, Celton J-M, Linsmith G, Becker C, Choisne N, Schijlen E, et al. High-quality *de novo* assembly of the apple genome and methylome dynamics of early fruit development. *Nat Genet*. 2017; 49: 1099–1106. <https://doi.org/10.1038/ng.3886> PMID: 28581499
51. Klagges C, Campoy JA, Quero-García J, Guzmán A, Mansur L, Gratacós E, et al. Construction and comparative analyses of highly dense linkage maps of two sweet cherry intra-specific progenies of commercial cultivars. *PLOS ONE*. 2013; 8: e54743. <https://doi.org/10.1371/journal.pone.0054743> PMID: 23382953
52. Chagné D, Kirk C, Whitworth C, Erasmuson S, Bicknell R, Sargent DJ, et al. Polyploid and aneuploid detection in apple using a single nucleotide polymorphism array. *Tree Genetics & Genomes*. 2015; 11: 94. <https://doi.org/10.1007/s11295-015-0920-8>
53. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet*. 2007; 81: 559–575. <https://doi.org/10.1086/519795> PMID: 17701901
54. Kalinowski ST, Taper ML, Marshall TC. Revising how the computer program CERVUS accommodates genotyping error increases success in paternity assignment. *Mol Ecol*. 2007; 16: 1099–1106. <https://doi.org/10.1111/j.1365-294X.2007.03089.x> PMID: 17305863
55. Young ND, Tanksley SD. Restriction fragment length polymorphism maps and the concept of graphical genotypes. *Theoret Appl Genetics*. 1989; 77: 95–101. <https://doi.org/10.1007/BF00292322> PMID: 24232480
56. Namjou B, Sestak AL, Armstrong DL, Zidovetzki R, Kelly JA, Jacob N, et al. High-density genotyping of *STAT4* reveals multiple haplotypic associations with systemic lupus erythematosus in different racial groups. *Arthritis Rheum*. 2009; 60: 1085–1095. <https://doi.org/10.1002/art.24387> PMID: 19333953
57. Jacob CO, Zhu J, Armstrong DL, Yan M, Han J, Zhou XJ, et al. Identification of *IRAK1* as a risk gene with critical role in the pathogenesis of systemic lupus erythematosus. *Proc Natl Acad Sci USA*. 2009; 106: 6256–6261. <https://doi.org/10.1073/pnas.0901181106> PMID: 19329491
58. Kumar S, Chagné D, Bink MCAM, Volz RK, Whitworth C, Carlisle C. Genomic selection for fruit quality traits in apple (*Malus × domestica* Borkh.). *PLOS ONE*. 2012; 7: e36674. <https://doi.org/10.1371/journal.pone.0036674> PMID: 22574211
59. Guo Y, He J, Zhao S, Wu H, Zhong X, Sheng Q, et al. Illumina human exome genotyping array clustering and quality control. *Nat Protoc*. 2014; 9: 2643–2662. <https://doi.org/10.1038/nprot.2014.174> PMID: 25321409
60. Vanderzande S, Micheletti D, Troggio M, Davey MW, Keulemans J. Genetic diversity, population structure, and linkage disequilibrium of elite and local apple accessions from Belgium using the IRSC array. *Tree Genetics & Genomes*. 2017; 13: 125. <https://doi.org/10.1007/s11295-017-1206-0>
61. da Silva Linge C, Antanaviciute L, Abdelghafar A, Arús P, Bassi D, Rossini L, et al. High-density multi-population consensus genetic linkage map for peach. *PLOS ONE*. 2018; 13: e0207724. <https://doi.org/10.1371/journal.pone.0207724> PMID: 30462743
62. Vanderzande S, Zheng P, Cai L, Iezzoni A, Main D, Peace C. Development and initial assessment of the 9K SNP addition to the sweet and sour cherry genome-wide SNP array. San Diego, CA, USA; 2019.
63. Okie WR. Handbook of peach and nectarine varieties: performance in the southeastern United States and index of names. U.S. Dept. of Agriculture, Agricultural Research Service; 1998.

64. Scorza R, Sherman W. Peaches. In: Janick J, Moore J, editors. Fruit Breeding. New York: Wiley; 1996. pp. 325–440.
65. Allard A, Bink MCAM, Martinez S, Kelner J-J, Legave, di Guardo M, et al. Detecting QTLs and putative candidate genes involved in budbreak and flowering time in an apple multiparental population. *J Exp Bot*. 2016; 67: 2875–2888. <https://doi.org/10.1093/jxb/erw130> PMID: 27034326
66. Di Guardo M, Bink MCAM, Guerra W, Letschka, Lozano L, Busatto N, et al. Deciphering the genetic control of fruit texture in apple by multiple family-based analysis and genome-wide association. *J Exp Bot*. 2017; 68: 1451–1466. <https://doi.org/10.1093/jxb/erx017> PMID: 28338805
67. Durand J-B, Allard A, Guitton B, van de Weg E, Bink MCAM, Costes E. Predicting flowering behavior and exploring its genetic determinism in an apple multi-family population based on statistical indices and simplified phenotyping. *Front Plant Sci*. 2017; 8: 858. <https://doi.org/10.3389/fpls.2017.00858> PMID: 28638387
68. Verma S, Evans K, Guan Y, Luby JJ, Rosyara UR, Howard NP, et al. Two large-effect QTLs, *Ma* and *Ma3*, determine genetic potential for acidity in apple fruit: Breeding insights from a multi-family study. *Tree Genetics & Genomes*. 2019; 15: 18.
69. Hernández Mora JR, Micheletti D, Bink M, Van de Weg E, Cantín C, Nazzicari N, et al. Integrated QTL detection for key breeding traits in multiple peach progenies. *BMC Genomics*. 2017; 18: 404. <https://doi.org/10.1186/s12864-017-3783-6> PMID: 28583082
70. Roach JA, Verma S, Peres NA, Jamieson AR, van de Weg WE, Bink MCAM, et al. FaRXf1: a locus conferring resistance to angular leaf spot caused by *Xanthomonas fragariae* in octoploid strawberry. *Theor Appl Genet*. 2016; 129: 1191–1201. <https://doi.org/10.1007/s00122-016-2695-1> PMID: 26910360
71. Mangandi J, Verma S, Osorio L, Peres NA, e E van de, Whitaker VM. Pedigree-Based Analysis in a multiparental population of octoploid strawberry reveals QTL alleles conferring resistance to *Phytophthora cactorum*. *G3: Genes, Genomes, Genetics*. 2017; 7: 1707–1719. <https://doi.org/10.1534/g3.117.042119> PMID: 28592652
72. Anciro A, Mangandi J, Verma S, Peres N, Whitaker VM, Lee S. FaRCg1: a quantitative trait locus conferring resistance to *Colletotrichum crown rot* caused by *Colletotrichum gloeosporioides* in octoploid strawberry. *Theor Appl Genet*. 2018; 131: 2167–2177. <https://doi.org/10.1007/s00122-018-3145-z> PMID: 30032317