



# **stuart: an R package for the curation of SNP genotypes from experimental crosses**

Marie Bourdon, Xavier Montagutelli

## **► To cite this version:**

Marie Bourdon, Xavier Montagutelli. stuart: an R package for the curation of SNP genotypes from experimental crosses. G3, 2022, 12, pp.jkac219. 10.1093/g3journal/jkac219 . pasteur-03854881

**HAL Id: pasteur-03854881**

**<https://pasteur.hal.science/pasteur-03854881>**

Submitted on 16 Nov 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# stuart: an R package for the curation of SNP genotypes from experimental crosses

Marie Bourdon  , Xavier Montagutelli  \*

Mouse Genetics Laboratory, Institut Pasteur, Université Paris Cité, F-75015 Paris, France

\*Corresponding author: Mouse Genetics Laboratory, Institut Pasteur, Université Paris Cité, F-75015 Paris, France. Email: [xavier.montagutelli@pasteur.fr](mailto:xavier.montagutelli@pasteur.fr)

## Abstract

Genetic mapping in 2-generation crosses requires genotyping, usually performed with single nucleotide polymorphism markers arrays which provide high-density genetic information. However, genetic analysis on raw genotypes can lead to spurious or unreliable results due to defective single nucleotide polymorphism assays or wrong genotype interpretation. Here, we introduce stuart, an open-source R package, which analyzes raw genotyping data to filter single nucleotide polymorphism markers based on informativeness, Mendelian inheritance pattern, and consistency with parental genotypes. The functions of this package provide a curation pipeline and formatting adequate for genetic analysis with the R/qtl package. stuart is available with detailed documentation from <https://gitlab.pasteur.fr/mouselab/stuart/>.

**Keywords:** R package; genetic analysis; SNP genotypes

## Introduction

Genetic mapping of Mendelian or quantitative traits in inbred strains is classically achieved in 2-generation crosses such as intercrosses (F2) and backcrosses (N2), in which the inheritance of the trait is compared with the genotypes at multiple genetic markers encompassing the genome map. Variations of a quantitative trait are controlled by one or more quantitative trait loci (QTLs). A QTL is defined as a marker at which individuals carrying different genotypes show different average trait values. QTL mapping searches for QTLs by testing the association between trait values and genotypes at markers spanning the genome map. The statistical significance of the association is expressed as logarithm of the odds (LOD) score which is calculated for each genotyped marker and, at intermediates positions, for pseudo-markers created by interval mapping, generating an LOD score curve (Broman 2001). The curve peaks at regions potentially associated with the trait. These peaks are called QTLs if they reach predefined statistical thresholds established either from general statistical models (Lander and Kruglyak 1995) or by permutation tests performed on the cross data. For each permutation, phenotypes are shuffled between individuals to break real associations, and LOD scores are calculated to identify peaks, which are all false positives. The distribution of the peak LOD scores over a large number (>1,000) of permutations provides statistical thresholds: if a LOD score of 3.8 or higher is observed in 5% of the permutations, this value will be taken as the  $P = 0.05$  threshold (Doerge and Churchill 1996). QTL mapping on F2s and N2s can be conducted with R packages such as R/qtl (Broman et al. 2003) and R/qtl2 (Broman et al. 2019).

With genome sequencing, single nucleotide polymorphisms (SNPs) have become the standard across species for their very

high frequency, low cost, and high-throughput analysis using various genotyping platforms. In mice, several generations of Mouse Universal Genotyping Arrays (MUGA) have been developed, the most recent being GigaMUGA (143k SNPs; Morgan et al. 2015) and MiniMUGA (10.8k SNPs; Sigmon et al. 2020). GigaMUGA provides high-density coverage for the fine characterization of inbred strains or outbred populations such as the Diversity Outbred (Svenson et al. 2012), while the modest number of SNPs in MiniMUGA is largely sufficient to genotype intercrossed or backcrossed individuals. However, SNP reliability is affected by the performance of genotyping platforms and polymorphism between and within inbred strains. Spurious or unreliable mapping outputs can result from defective SNP assays or wrong genotype interpretations. Therefore, raw data obtained from genotyping services must be curated before performing genetic analyses.

Several tools exist for quality control of SNP genotyping arrays, including Illumina's GenomeStudio. R packages such as argyle (Morgan 2015) analyze hybridization intensity signals from MUGA arrays. The simple genetic structure of 2-generation crosses provides specific and efficient means for identifying spurious genotyping data, such as consistency with parental genotypes and expected Mendelian proportions. The R/qtl package includes functions to build genetic maps and check for genotype consistency (<https://rqt.org/tutorials/geneticmaps.pdf>). However, this control is performed once genotypes have been imported and involves multiple steps of manual curation. To provide a more automated process of data curation before genetic analysis, we have developed stuart, an R package that implements a pipeline for automatic filtering and curation of SNP genotyping data from 2-generation crosses based on simple rules. This package formats raw SNP allele calls from Illumina files into genotypes ready for importation in R/qtl. Using 3 intercross datasets, we illustrate the consequences of inconsistent

Received: July 5, 2022. Accepted: August 19, 2022

© The Author(s) 2022. Published by Oxford University Press on behalf of Genetics Society of America.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

genotypes on the estimated marker map and QTL mapping, and how the curation achieved by each function in *stuart* leads to trustworthy results.

## Materials and methods

*stuart* is a tidyverse (Wickham et al. 2019) based R package requiring R version 3.5.0 or later. Its open source is available on Institut Pasteur's GitLab: <https://gitlab.pasteur.fr/mouselab/stuart/> and can be installed with devtools (Wickham et al. 2021). *stuart*'s vignette provides detailed descriptions of data import and of each function. *stuart* imports SNP allele calls from MUGA Illumina platform or other sources using the same file format. The central object of *stuart* is the marker table which summarizes for each marker, the alleles found in the population, the number of individuals of each genotype and the exclusion status resulting from the curation steps. *stuart* exports curated data to an R/qtl compatible format. The SNP annotation file used was downloaded from [https://raw.githubusercontent.com/kbroman/MUGAarrays/master/UWisc/mini\\_uwisc\\_v2.csv](https://raw.githubusercontent.com/kbroman/MUGAarrays/master/UWisc/mini_uwisc_v2.csv) (last accessed August 29, 2022).

Three datasets were used to test the package. This article presents the results from 176 (CC001/Unc X C57BL/6J-*Ifnar1* KO) F2 mice (dataset 1). The analysis of 2 other data sets, 94 (C57BL/6J-*Ifnar1* KO X 129S2/SvPas-*Ifnar1* KO) F2 mice (dataset 2) and 89 (C57BL/6NCrl X CC021/Unc) F2 mice (dataset 3) is presented as [Supplementary data](#). Quantitative traits were studied in the 3 F2s. Phenotype distributions are presented in [Supplementary Fig. 1](#). Genotyping was performed by Neogen (Auchincruive, Scotland) with MiniMUGA on DNA prepared from tail biopsies using standard phenol-chloroform extraction. Genotype call rate was 0.927, 0.931, and 0.948 for dataset 1, dataset 2, and dataset 3, respectively. QTL mapping was performed using R/qtl. Statistical significance of phenotype-genotype association was computed by data permutation (Doerge and Churchill 1996), which provides genome-wide thresholds accounting for multiple testing. The following thresholds were used, as commonly accepted (Members of the Complex Trait Consortium 2003):  $P = 0.05$  for significant association,  $P = 0.1$  and  $P = 0.63$  for the suggestive association. All figures were designed with ggplot2 (Wickham 2016) or R/qtl.

## Results and discussion

### Consequences of inconsistent genotypes

SNP data delivered by the Illumina platform are base alleles that need to be translated into genotypes for genetic analysis. From our experience on multiple 2-generation crosses, we identified several types of genotype inconsistencies that were responsible for distorted marker maps and spurious QTL mapping results. Recombination fraction (RF), which measures the genetic distance between 2 markers, is estimated in a cross by analyzing the proportion of recombinants between adjacent markers in all individuals. The map of markers calculated from the cross data should be consistent with their known positions. The R/qtl `est-map()` and `plotMap()` functions produce a graphical comparison of the 2 maps ([Fig. 1a](#) and [Supplementary Fig. 2, a and b](#)). For each chromosome, the known position of each marker provided in the annotation file (left) is connected with the estimated position (right) based on observed RF. With minimally curated genotypes (exclusion of nonpolymorphic markers and markers with over 50% missing genotypes), large RF was found in many instances between closely linked markers, resulting in fan-like patterns. To further describe these distortions, we computed the

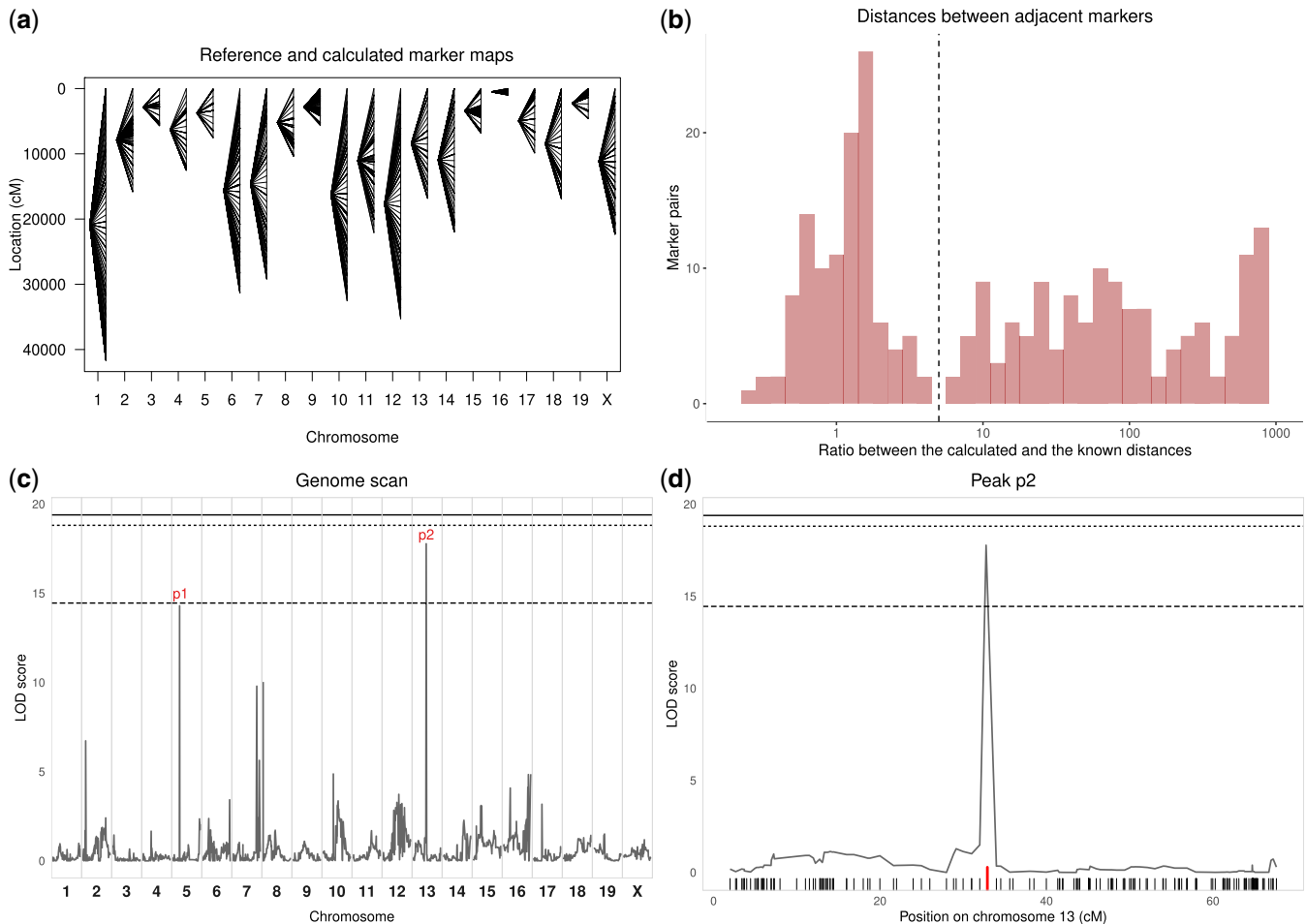
distribution of the ratio between the calculated and the known genetic distances between adjacent markers ([Fig. 1b](#) and [Supplementary Fig. 2, c and d](#); to avoid exaggerated ratios, we considered only markers with a known distance of 1 cM or more). This analysis revealed 2 groups of markers. In dataset 1, for 43% of them, the ratio was below 5 and followed a Gaussian distribution with mean = 1.31 and SD = 0.77. The other markers (57%) showed a ratio between 5 and 981.87 ([Fig. 1b](#)) which necessarily results from incorrect genotypes, as only a few individuals should show recombination between adjacent markers. On chromosome 1, while the known marker positions spanned ~100 cM, the cumulated genetic distance estimated from observed RF was ~40,000 cM. As QTL mapping relies on coherent genotypes at a series of markers encompassing a genetic interval, problematic genotypes at a given marker will perturb the analysis and, in some cases, may result in peaks of the LOD score curve in the absence of true association (Cheung et al. 2014). Such false positives increase significance thresholds calculated by data permutation.

These 2 consequences of genotyping inconsistencies are illustrated in [Fig. 1c](#) and [Supplementary Fig. 3, a and b](#) which were obtained using the R/qtl `scanone()` function on a quantitative trait from the uncured F2 datasets. For dataset 1, the  $P = 0.05$  significance threshold was estimated at 19.4 ([Fig. 1c](#)), while it usually ranges between 3.3 and 4.3 depending on the inheritance model for crosses of this type and size (Lander and Kruglyak 1995). Several peaks were detected although none reached  $P = 0.05$  significance. Moreover, their narrow profile was highly unexpected in F2 crosses. Indeed, these peaks involved only 1–3 markers, and the LOD score curve felt abruptly between these and adjacent markers on both sides ([Fig. 1d](#)), while genetic linkage between closely linked markers should result in progressive decrease of the LOD score curve on both sides of a peak (Guénet et al. 2015). Among the 3 datasets, we identified 4 narrow peaks reaching suggestive significance level ( $P < 0.63$ ): 2 were located at a marker with non-Mendelian allelic proportions and 2 were located at 1–3 pseudomarkers adjacent to a marker with non-Mendelian proportions ([Supplementary Fig. 3, c, d and e, f](#), respectively). We identified 5 other narrow peaks (LOD score between 6.72 and 10.03) out of which 4 resulted from the same situations as above and one was located on a pseudomarker and a marker with non-Mendelian proportions.

Inconsistent marker maps may also originate from the wrong assignment of markers to their chromosome and position provided to the mapping program. Indeed, R/qtl developer K. Broman identified errors in MUGA arrays annotation files affecting marker positions, probe sequences mapping to several locations, and unmappable markers. We recommend using K. Broman's corrected annotation files available on GitHub. The conversion of SNP alleles (A, C, T, G) observed in second-generation individuals (SGIs) to genotypes encoded according to the parental alleles may also create genotype errors. Reference SNP alleles established for many mouse strains may be used to infer the SGI genotypes. However, we recommend genotyping individuals of the parental strains used in the cross since they could differ from the reference panel. In our example dataset, the 2 parental strains used in the cross showed allelic differences with their reference panel counterpart at 200 markers.

### Data control and curation performed in *stuart*

Although each of *stuart*'s functions can be called independently, we present a logical analysis workflow appropriate for 2-generation crosses. [Table 1](#) summarizes the data curation and filtering



**Fig. 1.** Analysis of the dataset 1 illustrating the consequences of genotyping errors and inconsistencies on QTL mapping. Nonpolymorphic markers and markers with more than 50% missing genotypes were excluded to avoid excessive calculation time. a) Comparison of the known marker map (left) and the genetic map estimated from observed RF (right), as calculated by `est.map()` and represented by `plotMap()` functions of R/qtl. Lines connect the positions of each marker in the 2 maps. The estimated map is considerably expanded because of multiple genotype inconsistencies. b) Distribution of the ratio between estimated and known distances between adjacent markers. Markers with known and calculated distances below 1 cM were removed as they may lead to extremely small or large ratios. The expansion of the estimated map leads to a distribution tail of high ratios. The y-axis is in logarithmic scale. Fifty-seven percent of markers have a ratio above 5 (dashed line). c) Output of the scanone function of R/qtl showing the identification of narrow LOD score peaks. Genome-wide significance thresholds computed by data permutation are shown as plain ( $P = 0.05$ ), dotted ( $P = 0.1$ ), and dashed ( $P = 0.63$ ) lines. d) Magnification of the scanone plot restricted to chromosome 13 (peak p2). The LOD score peak is located on one marker (red tick) distant by 1.728 and 1.24 cM from the proximal and distal markers, respectively, on the known marker map, but by 1,001.582 and 1,001.506 cM based on calculated RF.

performed by each function, and the number of markers of dataset 1 retained after each step.

### Data importation

Genetic mapping requires both genotype and phenotype data. Required formats and instructions are detailed in the vignette (see example of phenotype data in [Supplementary Table 1](#)). Parental strains' genotyping data can be loaded from the same genotyping results as the SGI, from a previous genotyping file or from a reference file. Annotation data from K. Broman can be imported directly from GitHub. The `geno_strains()` function formats parental genotypes from a 2-allele encoding in Illumina format into a single letter encoding, and merges these data with the annotation table into a table with parental allele and marker positions.

### Consistency between parents and SGI alleles and genotypes

Several generations of MUGA arrays have been developed (Mega, Giga, Mini), each with successive versions differing by multiple

SNP markers. If parental and SGI data were produced on different versions, the marker lists must be compared to retain only common SNPs. This is achieved by the `mark_match()` function.

Converting alleles into genotypes requires that SGI segregate for the 2 parental alleles, and that each allele is found only in one parent. The aim of the `mark_allele()` function is to control consistency of allele's origin at multiple levels.

First, this function excludes markers with missing data in both parents. If allele data are available for only one parent and this allele is also found in SGI, the other allele present in SGI will be assigned to the parent with missing allele. However, this imputation is not error-free since we have observed, in rare occasions, markers which alleles were identical in the parental strains but were polymorphic in the SGI (Table 2 for such SNPs in dataset 1). This situation may occur when the parental strains used in the cross have diverged from those of the reference panel, or if one parent is heterozygous. Such markers will be excluded by the `mark_allele()` function but they could escape detection if allele information was missing in one parent.

**Table 1.** stuart analysis pipeline and application to dataset 1.

Steps	Function	Excluded markers	Number of markers retained
1. Import SGI alleles from MUGA arrays	read.table()/read_tsv()	–	11,125
2. Add data from parental strain			
Genotyped with SGI: make consensus	geno_strains()	–	–
Imported from another dataset: import and make consensus	read.table()/read_tsv(), geno_strains	–	–
Imported from reference	read.table()/other readr function depending on the format	–	–
3. Filter on allele consistency between parents and SGI			
Same set of markers between parents and SGI	mark_match()	Not present in both parents and SGI	11,125
Alleles consistent between parents and SGI	mark_allele	Missing alleles in both parents	10,375
		Not polymorphic in parents but polymorphic in SGI	
		Different alleles in parents and SGI	
		In backcrosses: homozygotes for the wrong allele	
		Optional: one parent missing or heterozygous	
4. Exclude markers with high proportion of missing genotypes	mark_na()	>50% of missing genotypes by default	9,918
5. Exclude nonpolymorphic markers in SGI	mark_poly()	Nonpolymorphic in SGI	2,738
6. Verify Mendelian proportions	mark_prop()	Departure from expected Mendelian segregation (proportion of each class or statistical threshold)	2,254
7. Verify RF between markers	est.map() followed by mark_estmap()	High RFs with adjacent markers	2,251

**Table 2.** Markers of dataset 1 non polymorphic between parental strains but polymorphic in SGI.

Marker	Allele parent 1	Allele parent 2	Allele SGI 1	Allele SGI 2
S6J017555686	C	C	T	C
S6J113080150	G	G	A	G
gJAX00038569	C	C	T	C
mUNC21540855	C	C	A	C
gUNC21555204	T	T	T	C
gUNC21596600	A	A	A	G

Adding the `parNH = FALSE` argument to the `mark_allele()` function will exclude markers missing one parental allele or for which one parent is heterozygous. However, while preventing rare errors, this option will also exclude a number of truly informative markers.

The `mark_allele()` function also discards markers at which parents and SGI carry different alleles, and, for backcrosses, markers for which some SGI are homozygous for the wrong allele.

### Nonpolymorphic markers

Genetic analysis requires polymorphic markers, i.e. for which parents carry different alleles which segregate in the SGI. The `mark_poly()` function excludes markers for which all genotyped SGI carry the same allele, which saves computation time.

### Missing genotypes

Reliable QTL mapping results depend on markers with medium to high rate of successful genotyping. Figure 2a shows markers distribution based on the proportion of missing genotypes. For over 95% of markers genotyping rate was above 50%. Genotyping

failures may result from poor-quality genotyping assay. The `mark_na()` function excludes such poorly genotyped markers.

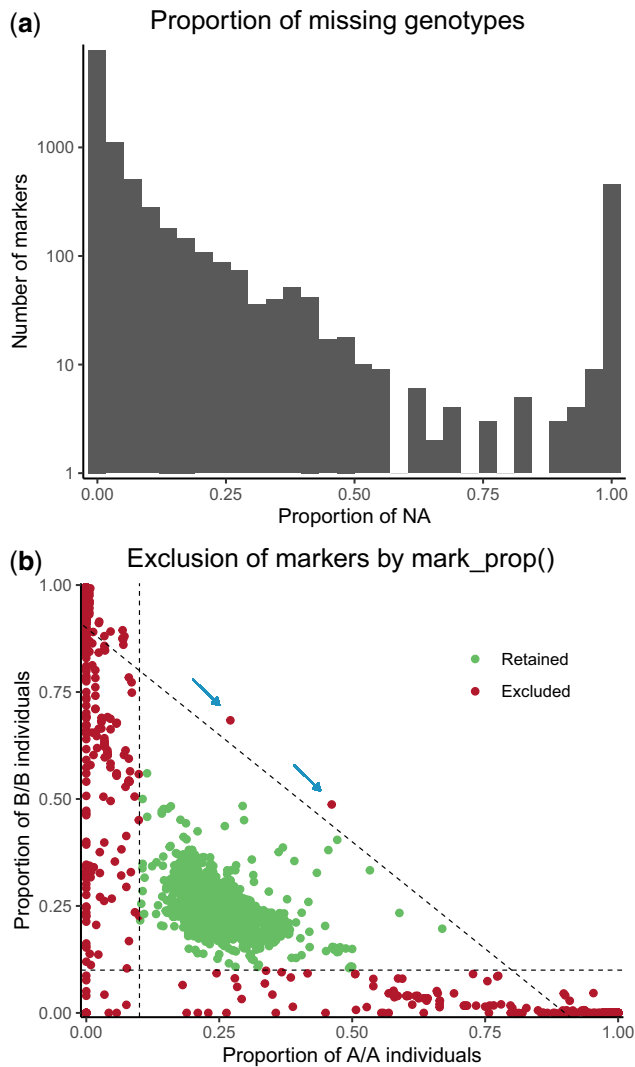
### Mendelian proportions

In 2-generation crosses between inbred strains, the proportions of the 2 or 3 classes of genotypes are predictable, i.e. for autosomes, 25% of each type of homozygotes and 50% of heterozygotes in an intercross, and 50% of homozygotes and 50% of heterozygotes in a backcross. Comparing the observed proportions with these expectations provides another criterion of filtering.

The `mark_prop()` function filters markers based either on a minimum proportion of each genotype or on the statistically significant departure from the expected proportions (Chi<sup>2</sup> test, with a P-value threshold). Figure 2b shows the exclusions of the autosomal markers depending on the proportion of each genotype. X chromosome genotypic proportions differ from autosomes, therefore, different arguments of `mark_prop()` function are used to filter X-linked markers for more precise curation.

### Filtering report and impact on QTL mapping results

At every step, the markers filtered out are annotated in a marker table which can be exported for further inspection. The last column of Table 1 shows the number of markers retained after each step in the example dataset 1. Most of the starting markers (7,180/11,125 = 65%) which were eventually removed by stuart's functions were removed by `mark_poly()` as nonpolymorphic, a ratio expected for crosses between 2 standard mouse inbred strains (Frazer et al. 2007). `mark_allele()` rejected 750 markers, `mark_na()` 457 and `mark_prop()` 484. Across the 3 datasets, we found 1,546 markers with either non-Mendelian proportions or allele inconsistencies between parental strains and SGIs. Overall, 619 of them were retained by stuart's filtering in at least one of the



**Fig. 2.** a) Distribution of the markers by their proportion of missing genotype (NA) in dataset 1. The y-axis is in logarithmic scale. 4.63% of markers have >50% missing genotypes. b) Exclusion of markers depending on genotypic proportions in dataset 1. Markers on X and Y chromosomes and mitochondrial DNA are not represented. The 2 axes represent the proportions of the 2 types of homozygous individuals in the intercross: AA and BB. Each dot represents a marker. Markers were excluded if the proportion of at least one of the 3 genotypes (AA, AB, and BB) was less than 10%, i.e. outside the triangle defined by the 3 dashed lines (AA = 0.1, BB = 0.1, and AA + BB = 0.9). Arrows point at 2 markers excluded due to a proportion of heterozygotes <10%.

other crosses, ruling out their misassignment to the genetic map. Out of the residual markers, 85 were removed from all datasets for another criterion than absence of polymorphism and were therefore considered as unreliable.

At this step, the dataset may still contain markers showing high RFs with adjacent markers either for a reason not tested by the current version of *stuart* or due to the parameters used in *mark\_na()* and *mark\_prop()* functions. These markers can be identified by calculating the estimated map using *R/qtl est.map()* and using *stuart*'s *mark\_estmap()* function which excludes markers presenting high RFs with adjacent markers. Over the 3 datasets, 9 markers were removed by *mark\_estmap()*. Five of them were retained in at least one other dataset, indicating the problem was dataset specific. Finally, for dataset 1, 2,251 markers passed all steps resulting in an average genetic interval between

adjacent markers lower than 2 cM, which is largely sufficient to perform QTL mapping (Darvasi *et al.* 1993). After curation, phenotype and genotype data are combined and exported in the *R/qtl* format using the *write\_rqtl()* function. The *qtl2convert* package (Broman 2021) converts this output into the adequate format required by the more recent *R/qtl2* package.

Figure 3a and Supplementary Fig. 4, a and b show the marker maps calculated after data curation with *stuart*. The known marker map and the estimated genetic map are consistent, with minimal expansions or contractions. Large ratios between the calculated and the known genetic distances between adjacent markers have been eliminated (Fig. 3b, Supplementary Fig. 4, c and d). QTL mapping analysis on curated dataset 1 is shown on Fig. 3c (to be compared with Fig. 1c; see Supplementary Fig. 5 for datasets 2 and 3). LOD thresholds are in the expected range for an F2, and the LOD score curve reveals broader peaks than in Fig. 1b, with progressive LOD score decrease on both sides of the peak marker. One significant and 3 suggestive QTLs were identified on chromosomes 12 (P-value = 0.037, Fig. 3d), 5 (P-value = 0.460), 10 (P-value = 0.157), and 15 (P-value = 0.244) which were not visible using noncurated data due to very high LOD score thresholds.

Being very simple to use and efficient at curating genotyping errors, *stuart* will facilitate the use of genotyping arrays for genetic mapping purposes in 2-generation crosses, bridging the gap between raw allele data produced by SNP platforms and genetic analysis software. Moreover, its functions can be used independently to analyze inbred strains genotypes. For example, *geno\_strain()* creates a genotype consensus between 2 or more individuals of the same strain suitable for further inspection, which can be useful when genotyping or regenotyping a strain of interest. Comparing genotyping results of an inbred strain after several generations of breeding with *mark\_allele()* will readily identify variants that have emerged or been selected over time. Likewise, this function will help identifying genetic variants between substrains.

## Web resources

The source code of the *stuart* package and the code used for the figures of this article are publicly available from <https://gitlab.pasteur.fr/mouselab/stuart/>.

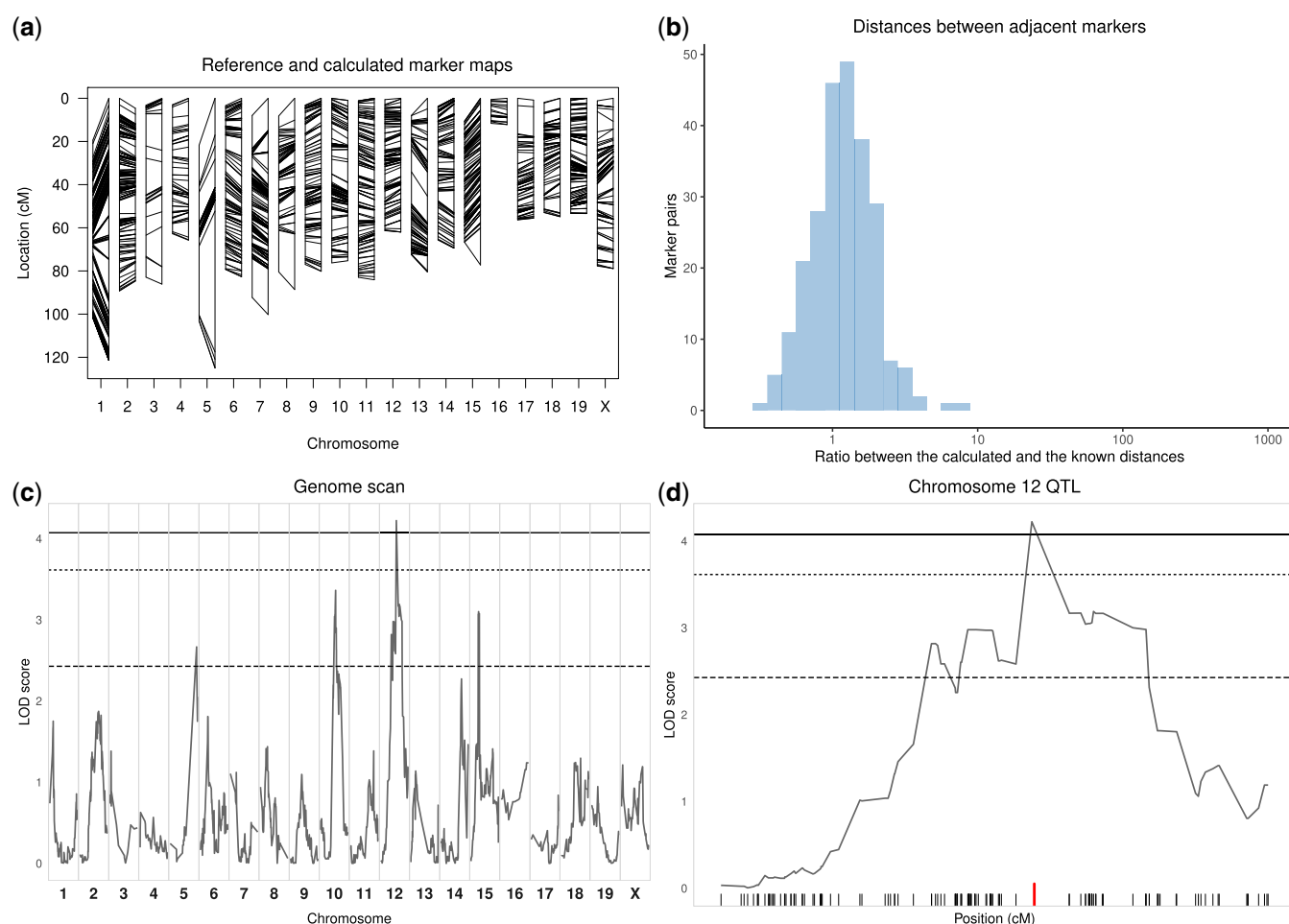
## Data availability

All datasets used as examples in this article are available from <https://gitlab.pasteur.fr/mouselab/stuart/>. Dataset 1 is included in the package and can be loaded once the package is loaded (see the vignette for details). The 2 other datasets are available from GitLab in the “article” directory in separate folders (i.e. “data2” and “data3”). Each folder contains the genotypes of the SGIs in file “geno\_dataX.csv,” the phenotypes of the SGIs in file “pheno\_dataX.csv,” the parental strains’ genotypes in file “parents\_dataX.csv” and the reference genotypes for the parental strains in file “ref\_gen\_dataX.csv.” Analysis of each cross is in each folder in an R markdown file (“dataX.Rmd”).

Supplemental material is available at G3 online.

## Acknowledgments

We thank Elise Jacquemet of the Pasteur Institute Bioinformatics and Biostatistics HUB for helping with the use of GitLab.



**Fig. 3.** Analysis of dataset 1 after curation of genotyping data by Stuart using the `mark_match()`, `mark_allele()`, `mark_na()`, `mark_poly()`, `mark_prop()`, and `mark_estmap()` functions. Refer to Fig. 1 for comparison with original data. a) The estimated marker map is now consistent with the known marker map. Despite some contraction or expansion of specific intervals, the genome length of the observed marker map for each chromosome is consistent with the known map (ratio between the calculated and the known length of the genome: 1.12). b) The distribution of the ratio between estimated and known distance between adjacent markers. Markers with known and calculated distances below 1 cM were removed as they may lead to extremely small or large ratios. Ratios are normally distributed with mean = 1.33 and SD = 0.81 showing consistency between the known and estimated maps. c) The LOD score curve shows several peaks, one of which is significant at  $P < 0.05$  (plain line, genome-wide significance computed by data permutation). Note that the significance thresholds are much lower than in Fig. 1c. None of the peaks shown in Fig. 1c were confirmed after data curation. Conversely, none of the peaks above  $P = 0.63$  (dashed line) found after data curation had been detected in Fig. 1c. d) The magnification of the QTL peak identified on chromosome 12, showing progressive decrease of the LOD score curve over a large genetic interval. The marker with the highest LOD score is identified with a thick (red) tick.

## Funding

This project was funded by the French Government's Investissement d'Avenir programme, Laboratoire d'Excellence "Integrative Biology of Emerging Infectious Diseases" (grant n°ANR-10-LABX-62-IBEID).

## Conflicts of interest

None declared.

## Literature cited

Broman KW. Review of statistical methods for QTL mapping in experimental crosses. *Lab Anim (NY)*. 2001;30(7):44–52.  
 Broman KW. `qtl2convert`: Convert Data among QTL Mapping Packages. 2021. [accessed 2022 May 15]. <https://CRAN.R-project.org/package=qtl2convert>.

Broman KW, Gatti DM, Simecek P, Furlotte NA, Prins P, Sen S, Yandell BS, Churchill GA. R/qtl2: software for mapping quantitative trait loci with high-dimensional data and multiparent populations. *Genetics*. 2019;211(2):495–502. doi:10.1534/genetics.118.301595.  
 Broman KW, Wu H, Sen S, Churchill GA. R/qtl: QTL mapping in experimental crosses. *Bioinformatics*. 2003;19(7):889–890. doi:10.1093/bioinformatics/btg112.  
 Cheung CYK, Thompson EA, Wijsman EM. Detection of Mendelian consistent genotyping errors in pedigrees: detection of genotyping errors. *Genet Epidemiol*. 2014;38(4):291–299. doi:10.1002/gepi.21806.  
 Darvasi A, Weinreb A, Minke V, Weller JI, Soller M. Detecting marker-QTL linkage and estimating QTL gene effect and map location using a saturated genetic map. *Genetics*. 1993;134(3):943–951. doi:10.1093/genetics/134.3.943.  
 Doerge RW, Churchill GA. Permutation tests for multiple loci affecting a quantitative character. *Genetics*. 1996;142(1):285–294. doi:10.1093/genetics/142.1.285.  
 Frazer KA, Eskin E, Kang HM, Bogue MA, Hinds DA, Beilharz EJ, Gupta RV, Montgomery J, Morenzoni MM, Nilsen GB, et al. A sequence-

- based variation map of 8.27 million SNPs in inbred mouse strains. *Nature*. 2007;448(7157):1050–1053. doi:[10.1038/nature06067](https://doi.org/10.1038/nature06067).
- Guénet JL, Benavides F, Panthier J-J, Montagutelli X. 2015. Genetics of the Mouse. [accessed 2022 Jun 24]. <https://link.springer.com/book/10.1007/978-3-662-44287-6>.
- Lander E, Kruglyak L. Genetic dissection of complex traits: guidelines for interpreting and reporting linkage results. *Nat Genet*. 1995; 11(3):241–247. doi:[10.1038/ng1195-241](https://doi.org/10.1038/ng1195-241).
- Members of the Complex Trait Consortium. The nature and identification of quantitative trait loci: a community's view. *Nat Rev Genet*. 2003;4(11):911–916. doi:[10.1038/nrg1206](https://doi.org/10.1038/nrg1206).
- Morgan AP. argyle: an R package for analysis of illumina genotyping arrays. *G3 (Bethesda)*. 2015;6(2):281–286. doi:[10.1534/g3.115.023739](https://doi.org/10.1534/g3.115.023739).
- Morgan AP, Fu C-P, Kao C-Y, Welsh CE, Didion JP, Yadgary L, Hyacinth L, Ferris MT, Bell TA, Miller DR, et al. The mouse universal genotyping array: from substrains to subspecies. *G3 (Bethesda)*. 2015;6(2):263–279. doi:[10.1534/g3.115.022087](https://doi.org/10.1534/g3.115.022087).
- Sigmon JS, Blanchard MW, Baric RS, Bell TA, Brennan J, Brockmann GA, Burks AW, Calabrese JM, Caron KM, Cheney RE, et al. Content and performance of the MiniMUGA genotyping array: a new tool to improve rigor and reproducibility in mouse research. *Genetics*. 2020;216(4):905–930. doi:[10.1534/genetics.120.303596](https://doi.org/10.1534/genetics.120.303596).
- Svenson KL, Gatti DM, Valdar W, Welsh CE, Cheng R, Chesler EJ, Palmer AA, McMillan L, Churchill GA. High-resolution genetic mapping using the mouse diversity outbred population. *Genetics*. 2012;190(2):437–447. doi:[10.1534/genetics.111.132597](https://doi.org/10.1534/genetics.111.132597).
- Wickham H. ggplot2: elegant Graphics for Data Analysis. 2016. [accessed 2022 May 15]. <https://ggplot2.tidyverse.org>.
- Wickham H, Averick M, Bryan J, Chang W, McGowan L, François R, Grolemund G, Hayes A, Henry L, Hester J, et al. Welcome to the tidyverse. *J Open Source Softw*. 2019;4(43):1686. doi:[10.21105/joss.01686](https://doi.org/10.21105/joss.01686).
- Wickham H, Hester J, Chang W, Bryan J. devtools: Tools to Make Developing R Packages Easier. 2021. [accessed 2022 May 15]. <https://devtools.r-lib.org>.

Communicating editor: F. P.-M. de Villena