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

2

3    **Authors**

4    **Marie Bourdon[1], Xavier Montagutelli[1]**

5

6    1 Institut Pasteur, Université Paris Cité, Mouse Genetics Laboratory, Paris, France

7    Correspondance : xavier.montagutelli@pasteur.fr

8

9    ORCID numbers

10    Marie Bourdon: 0000-0002-6184-1888

11    Xavier Montagutelli: 0000-0002-9372-5398

12

13    **Keywords**

14    R-package, genetic analysis, SNP genotypes

15

16    **Abstract**

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

25    **Introduction**

26    Genetic mapping of Mendelian or quantitative traits in inbred strains is classically achieved in two-
27    generation crosses such as intercrosses (F2) and backcrosses (N2), in which the inheritance of the
28    trait is compared with the genotypes at multiple genetic markers encompassing the genome map.
29    Variations of a quantitative trait are controlled by one or more quantitative trait loci (QTL). A QTL is
30    defined as a marker at which individuals carrying different genotypes show different average trait
31    values. QTL mapping searches for QTLs by testing association between trait values and genotypes at
32    markers spanning the genome map. The statistical significance of the association is expressed as
33    logarithm of the odds (LOD) score which is calculated for each genotyped marker and, at
34    intermediates positions, for pseudomarkers created by interval mapping, generating a LOD score
35    curve (Broman 2001). The curve peaks at regions potentially associated with the trait. These peaks
36    are called QTLs if they reach predefined statistical thresholds established either from general
37    statistical models (Lander and Kruglyak 1995) or by permutation tests performed on the cross data.

38  For each permutation, phenotypes are shuffled between individuals to break real associations, and
39  LOD scores are calculated to identify peaks, which are all false positives. The distribution of the peak
40  LOD scores over a large number (>1000) of permutations provides statistical thresholds: if a LOD
41  score of 3.8 or higher is observed in 5% of the permutations, this value will be taken as the p=0.05
42  threshold (Doerge and Churchill 1996). QTL mapping on F2s and N2s can be conducted with R
43  packages such as R/qtl (Broman *et al.* 2003) and R/qtl2 (Broman *et al.* 2019).

44

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

57  Several tools exist for quality control of SNP genotyping arrays, including Illumina's GenomeStudio. R
58  packages such as argyle (Morgan 2016) analyze hybridization intensity signals from MUGA arrays.
59  The simple genetic structure of two-generation crosses provides specific and efficient means for
60  identifying spurious genotyping data, such as consistency with parental genotypes and expected
61  Mendelian proportions. The R/qtl package includes functions to build genetic maps and check for
62  genotype   consistency   (https://rqtl.org/tutorials/geneticmaps.pdf).   However,   this   control   is
63  performed once genotypes have been imported and involves multiple steps of manual curation. To
64  provide a more automated process of data curation before genetic analysis, we have developed
65  stuart, an R package which implements a pipeline for automatic filtering and curation of SNP
66  genotyping data from two-generation crosses based on simple rules. This package formats raw SNP
67  allele calls from Illumina files into genotypes ready for importation in R/qtl. Using three intercross
68  datasets, we illustrate the consequences of inconsistent genotypes on the estimated marker map
69  and QTL mapping, and how the curation achieved by each function in stuart leads to trustable
70  results.

71

## **Materials and Methods**

73  stuart is a tidyverse (Wickham *et al.* 2019) based R package requiring R version 3.5.0 or later. Its
74  open source is available on Institut Pasteur's GitLab: https://gitlab.pasteur.fr/mouselab/stuart/ and
75  can be installed with devtools (Wickham *et al.* 2021). stuart's vignette provides detailed descriptions
76  of data import and of each function.

77  stuart imports SNP allele calls from MUGA Illumina platform or other sources using the same file
78  format. The central object of stuart is the marker table which summarizes for each marker the alleles
79  found in the population, the number of individuals of each genotype and the exclusion status
80  resulting from the curation steps. stuart exports curated data to an R/qtl compatible format. The

81  SNP annotation file used was downloaded from
82  https://raw.githubusercontent.com/kbroman/MUGAarrays/master/UWisc/mini_uwisc_v2.csv.

83  Three datasets were used to test the package. This article presents the results from 176 (CC001/Unc
84  X C57BL/6J-*Ifnar1* KO) F2 mice (dataset 1). The analysis of two other data sets, 94(C57BL/6J-*Ifnar1*
85  KO X 129S2/SvPas-*Ifnar1* KO) F2 mice (dataset 2) and 89 (C57BL/6NCrl X CC021/Unc) F2 mice
86  (dataset 3) is presented as supplementary data. Quantitative traits were studied in the three F2s.
87  Phenotype distributions are presented in Supplementary Figure 1. Genotyping was performed by
88  Neogen (Auchincruive, Scotland) with MiniMUGA on DNA prepared from tail biopsies using standard
89  phenol-chloroform extraction. Genotype call rate was 0.927, 0.931 and 0.948 for dataset 1, dataset 2
90  and dataset 3, respectively. QTL mapping was performed using R/qtl. The following thresholds were
91  used, as commonly accepted (Members of the Complex Trait Consortium 2003): p=0.05 for
92  significant association, p=0.1 and p=0.63 for suggestive association. All figures were designed with
93  ggplot2 (Wickham 2016) or R/qtl.

94

95  **Results and discussion**

96  Consequences of inconsistent genotypes

97  SNP data delivered by the Illumina platform are base alleles that need to be translated into
98  genotypes for genetic analysis. From our experience on multiple two-generation crosses, we
99  identified several types of genotype inconsistencies which were responsible for distorted marker
100  maps and spurious QTL mapping results. Recombination fraction (RF), which measures the genetic
101  distance between two markers, is estimated in a cross by analyzing proportion of recombinants
102  between adjacent markers in all individuals. The map of markers calculated from the cross data
103  should be consistent with their known positions. The R/qtl est.map() and plotMap() functions
104  produce a graphical comparison of the two maps (Figure 1A and Supplementary Figure 2A-B). For
105  each chromosome, the known position of each marker provided in the annotation file (left) is
106  connected with the estimated position (right) based on observed RF. With minimally curated
107  genotypes (exclusion of non-polymorphic markers and markers with over 50% missing genotypes),
108  large RF were found in many instances between closely linked markers, resulting in fan-like patterns.
109  To further describe these distortions, we computed the distribution of the ratio between the
110  calculated and the known genetic distances between adjacent markers (Figure 1B and
111  Supplementary Figure 2C-D; to avoid exaggerated ratios, we considered only markers with a known
112  distance of 1cM or more). This analysis revealed two groups of markers. In dataset 1, for 43%
113  percent of them, the ratio was below 5 and followed a Gaussian distribution with mean = 1.31 and
114  sd = 0.77. The other markers (57%) showed a ratio between 5 and 981.87 (Figure 1B) which
115  necessarily results from incorrect genotypes, as only a few individuals should show recombination
116  between adjacent markers. On chromosome 1, while the known marker positions spanned ~100 cM,
117  the cumulated genetic distance estimated from observed RF was ~40,000cM. As QTL mapping relies
118  on coherent genotypes at a series of markers encompassing a genetic interval, problematic
119  genotypes at a given marker will perturb the analysis and, in some cases, may result in peaks of the
120  LOD score curve in the absence of true association (Cheung *et al.* 2014). Such false positives increase
121  significance thresholds calculated by data permutation.

122  These two consequences of genotyping inconsistencies are illustrated in Figure 1C and
123  Supplementary Figure 3A-B which were obtained using the R/qtl scanone() function on a
124  quantitative trait from the uncurated F2 datasets. For dataset 1, the p=0.05 significance threshold

125   was estimated at 19.4 (Figure 1C), while it usually ranges between 3.3 and 4.3 depending on the
126   inheritance model for crosses of this type and size (Lander and Kruglyak 1995). Several peaks were
127   detected although none reached p=0.05 significance. Moreover, their narrow profile was highly
128   unexpected in F2 crosses. Indeed, these peaks involved only one to three markers, and the LOD
129   score curve felt abruptly between these and adjacent markers on both sides (Figure 1D), while
130   genetic linkage between closely linked markers should result in progressive decrease of the LOD
131   score curve on both sides of a peak (Guénet *et al.* 2015). Among the three datasets, we identified
132   four narrow peaks reaching suggestive significance level (p<0.63): two were located at a marker with
133   non-Mendelian allelic proportions and two were located at one to three pseudomarkers adjacent to
134   a marker with non-Mendelian proportions (Supplementary Figure 3C-D and E-F, respectively). We
135   identified 5 other narrow peaks (LOD score between 6.72 and 10.03) out of which four resulted from
136   the same situations as above and one was located on a pseudomaker and a marker with non-
137   Mendelian proportions.

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

149

150   <u>Data control and curation performed in stuart</u>

151   Although each of stuart's functions can be called independently, we present a logical analysis
152   workflow appropriate for two-generation crosses. Table 1 summarizes the data curation and filtering
153   performed by each function, and the number of markers of dataset 1 retained after each step.

154   *Data importation*

155   Genetic mapping requires both genotype and phenotype data. Required formats and instructions
156   are detailed in the vignette (see example of phenotype data in Supplementary Table S1). Parental
157   strains' genotyping data can be loaded from the same genotyping results as the SGI, from a previous
158   genotyping file or from a reference file. Annotation data from K. Broman can be imported directly
159   from GitHub. The geno_strains() function formats parental genotypes from a two-allele encoding in
160   Illumina format into a single letter encoding, and merges this data with the annotation table into a
161   table with parental allele and marker positions.

162   *Consistency between parents and SGI alleles and genotypes*

163   Several generations of MUGA arrays have been developed (Mega, Giga, Mini), each with successive
164   versions differing by multiple SNP markers. If parental and SGI data were produced on different
165   versions, the marker lists must be compared to retain only common SNPs. This is achieved by the
166   mark_match() function.

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

170  First, this function excludes markers with missing data in both parents. If allele data is available for
171  only one parent and this allele is also found in SGI, the other allele present in SGI will be assigned to
172  the parent with missing allele. However, this imputation is not error-free since we have observed, in
173  rare occasions, markers which alleles were identical in the parental strains but were polymorphic in
174  the SGI (Table 2 for such SNPs in dataset 1). This situation may occur when the parental strains used
175  in the cross have diverged from those of the reference panel, or if one parent is heterozygous. Such
176  markers will be excluded by the mark_allele() function but they could escape detection if allele
177  information was missing in one parent. Adding the parNH=FALSE argument to the mark_allele()
178  function will exclude markers missing one parental allele or for which one parent is heterozygous.
179  However, while preventing rare errors, this option will also exclude a number of truly informative
180  markers.

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

*Non-polymorphic markers*

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

*Missing genotypes*

188  Reliable QTL mapping results depend on markers with medium to high rate of successful genotyping.
189  Figure 2A shows markers distribution based on the proportion of missing genotypes. For over 95% of
190  markers genotyping rate was above 50%. Genotyping failures may result from poor-quality
191  genotyping assay. The mark_na() function excludes such poorly-genotyped markers.

*Mendelian proportions*

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

197  The mark_prop() function filters markers based either on a minimum proportion of each genotype,
198  or on the statistically significant departure from the expected proportions (Chi2 test, with a p-value
199  threshold). Figure 2B shows the exclusions of the autosomal markers depending on the proportion
200  of each genotype. X chromosome genotypic proportions differ from autosomes, therefore, different
201  arguments of mark_prop() function are used to filter X-linked markers for more precise curation.

*Filtering report and impact on QTL mapping results*

203  At every step, the markers filtered out are annotated in a marker table which can be exported for
204  further inspection. The last column of Table 1 shows the number of markers retained after each step
205  in the example dataset 1. Most of the starting markers (7180/11125 = 65%) which were eventually
206  removed by stuart's functions were removed by mark_poly() as non-polymorphic, a ratio expected
207  for crosses between two standard mouse inbred strains (Frazer *et al.* 2007). mark_allele() rejected
208  750 markers, mark_na() 457 and mark_prop() 484. Across the three datasets, we found 1546

209 markers with either non-Mendelian proportions or allele inconsistencies between parental strains
210 and SGIs. Overall, 619 of them were retained by stuart's filtering in at least one of the other crosses,
211 ruling out their misassignment to the genetic map. Out of the residual markers, 85 were removed
212 from all datasets for another criterion than absence of polymorphism and were therefore
213 considered as unreliable.

214 At this step, the dataset may still contain markers showing high recombination fractions with
215 adjacent markers either for a reason not tested by the current version of stuart or due to the
216 parameters used in mark_na() and mark_prop() functions. These markers can be identified by
217 calculating the estimated map using R/qtl est.map()and using stuart's mark_estmap() function which
218 excludes markers presenting high recombination fractions with adjacent markers. Over the 3
219 datasets, 9 markers were removed by mark_estmap(). Five of them were retained in at least one
220 other dataset, indicating the problem was dataset specific. Finally, for dataset 1, 2251 markers
221 passed all steps resulting in an average genetic interval between adjacent markers lower than 2 cM,
222 which is largely sufficient to perform QTL mapping (Darvasi *et al.* 1993). After curation, phenotype
223 and genotype data are combined and exported in the R/qtl format using the write_rqtl() function.
224 The qtl2convert package (Broman 2021) converts this output into the adequate format required by
225 the more recent R/qtl2 package.

226 Figure 3A and Supplementary Figure 4A-B show the marker maps calculated after data curation with
227 stuart. The known marker map and the estimated genetic map are consistent, with minimal
228 expansions or contractions. Large ratios between the calculated and the known genetic distances
229 between adjacent markers have been eliminated (Figure 3B, Supplementary Figure 4C-D). QTL
230 mapping analysis on curated dataset 1 is shown on Figure 3C (to be compared with Figure 1C; see
231 Supplementary Figure 5 for datasets 2 and 3). LOD thresholds are in the expected range for an F2,
232 and the LOD score curve reveals broader peaks than in Figure 1B, with progressive LOD score
233 decrease on both sides of the peak marker. One significant and three suggestive QTLs were
234 identified on chromosomes 12 (p-value = 0.037, Figure 3D), 5 (p-value = 0.460), 10 (p-value = 0.157)
235 and 15 (p-value = 0.244) which were not visible using non-curated data due to very high LOD score
236 thresholds.

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

246

247 **Web resources**

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

250

251 **Data availability statement**

6

252    All datasets used as examples in this article are available from
253    https://gitlab.pasteur.fr/mouselab/stuart/. Dataset 1 is included in the package and can be loaded
254    once the package is loaded (see the vignette for details). The two other datasets are available from
255    GitLab in the "article" directory in separate folders (i.e. "data2" and "data3"). Each folder contains
256    the genotypes of the SGIs in file "geno_dataX.csv", the phenotypes of the SGIs in file
257    "pheno_dataX.csv", the parental strains' genotypes in file "parents_dataX.csv" and the reference
258    genotypes for the parental strains in file "ref_geno_dataX.csv". Analysis of each cross is in each
259    folder in an R markdown file ("dataX.Rmd").

260

**Acknowledgements**

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

264

**Conflict of interest**

266    The authors declare no conflicting interests.

267

**Funder information**

272

273 **Literature cited**

274 Broman, K. W., 2021 *qtl2convert: Convert Data among QTL Mapping Packages.*

275 Broman, K. W., 2001 Review of statistical methods for QTL mapping in experimental crosses. Lab

276       Anim (NY) 30: 44–52.

277 Broman, K. W., D. M. Gatti, P. Simecek, N. A. Furlotte, P. Prins *et al.*, 2019 R/qtl2: Software for

278       Mapping Quantitative Trait Loci with High-Dimensional Data and Multiparent Populations.

279       Genetics 211: 495–502.

280 Broman, K. W., H. Wu, S. Sen, and G. A. Churchill, 2003 R/qtl: QTL mapping in experimental crosses.

281       Bioinformatics 19: 889–890.

282 Cheung, C. Y. K., E. A. Thompson, and E. M. Wijsman, 2014 Detection of Mendelian Consistent

283       Genotyping Errors in Pedigrees: Detection of Genotyping Errors. Genet. Epidemiol. 38: 291–

284       299.

285 Darvasi, A., A. Weinreb, V. Minke, J. I. Weller, and M. Soller, 1993 Detecting marker-QTL linkage and

286       estimating QTL gene effect and map location using a saturated genetic map. Genetics 134:

287       943–951.

288 Doerge, R. W., and G. A. Churchill, 1996 Permutation Tests for Multiple Loci Affecting a Quantitative

289       Character. Genetics 142: 285–294.

290 Frazer, K. A., E. Eskin, H. M. Kang, M. A. Bogue, D. A. Hinds *et al.*, 2007 A sequence-based variation

291       map of 8.27 million SNPs in inbred mouse strains. Nature 448: 1050–1053.

292 Guénet, J. L., F. Benavides, J.-J. Panthier, and X. Montagutelli, 2015 *Genetics of the Mouse.*

293 Lander, E., and L. Kruglyak, 1995 Genetic dissection of complex traits: guidelines for interpreting and

294       reporting linkage results. Nat Genet 11: 241–247.

295 Members of the Complex Trait Consortium, 2003 The nature and identification of quantitative trait

296       loci: a community's view. Nat Rev Genet 4: 911–916.

297 Morgan, A. P., 2016 argyle: An R Package for Analysis of Illumina Genotyping Arrays. G3

298       Genes|Genomes|Genetics 6: 281–286.

299    Morgan, A. P., C.-P. Fu, C.-Y. Kao, C. E. Welsh, J. P. Didion *et al.*, 2015 The Mouse Universal

300          Genotyping Array: From Substrains to Subspecies. G3 (Bethesda) 6: 263–279.

301    Sigmon, J. S., M. W. Blanchard, R. S. Baric, T. A. Bell, J. Brennan *et al.*, 2020 Content and Performance

302          of the MiniMUGA Genotyping Array: A New Tool To Improve Rigor and Reproducibility in

303          Mouse Research. Genetics 216: 905–930.

304    Svenson, K. L., D. M. Gatti, W. Valdar, C. E. Welsh, R. Cheng *et al.*, 2012 High-Resolution Genetic

305          Mapping Using the Mouse Diversity Outbred Population. Genetics 190: 437–447.

306    Wickham, H., 2016 *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York.

307    Wickham, H., M. Averick, J. Bryan, W. Chang, L. McGowan *et al.*, 2019 Welcome to the Tidyverse.

308          JOSS 4: 1686.

309    Wickham, H., J. Hester, W. Chang, and J. Bryan, 2021 *devtools: Tools to Make Developing R Packages*

310          *Easier*.

311

312

313 **Table 1 – stuart analysis pipeline and application to dataset 1**

314

| Steps | Function | Excluded markers | Number of markers retained |
|---|---|---|---|
| 1. Import SGI alleles from MUGA arrays | read.table()/read_tsv() | - | 11125 |
| 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 | 11125 |
| Alleles consistent between parents and SGI | mark_allele | Missing alleles in both parents<br>Not polymorphic in parents but polymorphic in SGI<br>Different alleles in parents and SGI<br>In backcrosses: homozygotes for the wrong allele<br>Optional: one parent missing or heterozygous | 10375 |
| 4. Exclude markers with high proportion of missing genotypes | mark_na() | >50% of missing genotypes by default | 9918 |
| 5. Exclude non-polymorphic markers in SGI | mark_poly() | Non polymorphic in SGI | 2738 |
| 6. Verify mendelian proportions | mark_prop() | Departure from expected Mendelian segregation (proportion of each class or statistical threshold) | 2254 |
| 7. Verify recombination fraction between markers | est.map() followed by mark_estmap() | High recombination fractions with adjacent markers | 2251 |

315

10

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

317

| 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 |

318

319

320    **Figure legends**

321    **Figure 1**

322    Analysis of the dataset 1 illustrating the consequences of genotyping errors and inconsistencies on
323    QTL mapping. Non-polymorphic markers and markers with more than 50% missing genotypes were
324    excluded to avoid excessive calculation time. A: comparison of the known marker map (left) and the
325    genetic map estimated from observed RF (right), as calculated by est.map() and represented by
326    plotMap() functions of R/qtl. Lines connect the positions of each marker in the two maps. The
327    estimated map is considerably expanded because of multiple genotype inconsistencies. B:
328    distribution of the ratio between estimated and known distances between adjacent markers.
329    Markers with known and calculated distances below 1cM were removed as they may lead to
330    extremely small or large ratios. The expansion of the estimated map leads to a distribution tail of
331    high ratios. The y-axis is in logarithmic scale. 57% of markers have a ratio above 5 (dashed line). C:
332    output of the scanone function of R/qtl showing the identification of narrow LOD score peaks.
333    Significance thresholds are shown as plain (p=0.05), dotted (p=0.1) and dashed (p=0.63) lines. D:
334    magnification of the scanone plot restricted to chromosome 13 (peak p2). The LOD score peak is
335    located on one marker (red tick) distant by 1.728 cM and 1.24 cM from the proximal and distal
336    markers, respectively, on the known marker map, but by 1001.582 cM and 1001.506 cM based on
337    calculated RF.

338

339    **Figure 2**

340    A: Distribution of the markers by their proportion of missing genotype (NA) in dataset 1. The y-axis is
341    in logarithmic scale. 4.63% of markers have >50% missing genotypes. B: Exclusion of markers
342    depending on genotypic proportions in dataset 1. Markers on X and Y chromosomes and
343    mitochondrial DNA are not represented. The two axes represent the proportions of the two types of
344    homozygous individuals in the intercross: AA and BB. Each dot represents a marker. Markers were
345    excluded (red) if the proportion of at least one of the three genotypes (AA, AB and BB) was less than
346    10%, i.e. outside the triangle defined by the three dashed lines (AA=0.1, BB=0.1 and AA+BB=0.9).
347    Blue arrows point at two markers excluded due to a proportion of heterozygotes <10%.

348

349    **Figure 3**

350    Analysis of dataset 1 after curation of genotyping data by stuart using the mark_match(),
351    mark_allele(), mark_na(), mark_poly(), mark_prop(), and mark_estmap() functions. Refer to Figure 1
352    for comparison with original data. A: the estimated marker map is now consistent with the known
353    marker map. Despite some contraction or expansion of specific intervals, the genome length of the
354    observed marker map for each chromosome is consistent with the known map (ratio between the
355    calculated and the known length of the genome:1.12). B: distribution of the ratio between estimated
356    and known distance between adjacent markers. Markers with known and calculated distances below
357    1cM were removed as they may lead to extremely small or large ratios. Ratios are normally
358    distributed with mean=1.33 and sd=0.81 showing consistency between the known and estimated
359    maps. C: The LOD score curve shows several peaks, one of which is significant at P<0.05 (plain line).
360    Note that the significance thresholds are much lower than in Figure 1C. None of the peaks shown in
361    Figure1C were confirmed after data curation. Conversely, none of the peaks above P=0.63 (dashed

362    line) found after data curation had been detected in Figure1C. D: magnification of the QTL peak
363    identified on chromosome 12, showing progressive decrease of the LOD score curve over a large
364    genetic interval. The marker with the highest LOD score is identified with a red tick.

365

366 **Supplementary data**

367

368 **Supplementary Table 1**

369 Format of the phenotype data. First column: individual's number. Second column: individual's sex.
370 Following columns: traits and covariates (individual's age, phenotype)

371 **Supplementary Figure 1**

372 Distribution of the quantitative phenotypes analyzed in the three datasets.

373 **Supplementary Figure 2**

374 Analysis of datasets 2 and 3 illustrating the expansion of the estimated genetic maps. Non-
375 polymorphic markers and markers with more than 50% missing genotypes were excluded to avoid
376 excessive calculation time. A, B: comparison of the known marker map (left) and the genetic map
377 estimated from observed RF (right), as calculated by est.map() and represented by plotMap()
378 functions of R/qtl in dataset 2 (A) and dataset 3 (B). Lines connect the positions of each marker in
379 the two maps. The estimated map is considerably expanded because of multiple genotype
380 inconsistencies. C, D: distribution of the ratio between estimated and known distance between
381 adjacent markers. Markers with known and calculated distances below 1cM were removed as they
382 may lead to extremely small or large ratios. The expansion of the estimated map leads to a
383 distribution tail of high ratios. The y-axis is in logarithmic scale. 51% of the markers in the dataset 2
384 (C) and 16% of the markers in the dataset 3 (D) have a ratio above 5 (dashed line).

385 **Supplementary Figure 3**

386 Analysis of datasets 2 and 3 illustrating the identification of narrow LOD-score peaks. A, B: output of
387 the scanone() function of R/qtl in the datasets 2 (A) and 3 (B) showing the identification of two
388 narrow suggestive peaks. C: peak p1 from dataset 1 (see Figure 1C) is located on a single marker
389 (mUNC050096588, red tick) with non-Mendelian proportions. Peak p2 from dataset 1 shows the
390 same pattern. D: genotypes at mUNC050096588. HM1 and HM2: homozygotes; HT: heterozygotes;
391 NA: missing genotypes. E: peak p3 from dataset 2 is located on a pseudomarker adjacent to a
392 marker with non-Mendelian proportions (SNT111392585, red tick). Peak p4 from dataset 3 shows
393 the same pattern. F: genotypes at SNT111392585.

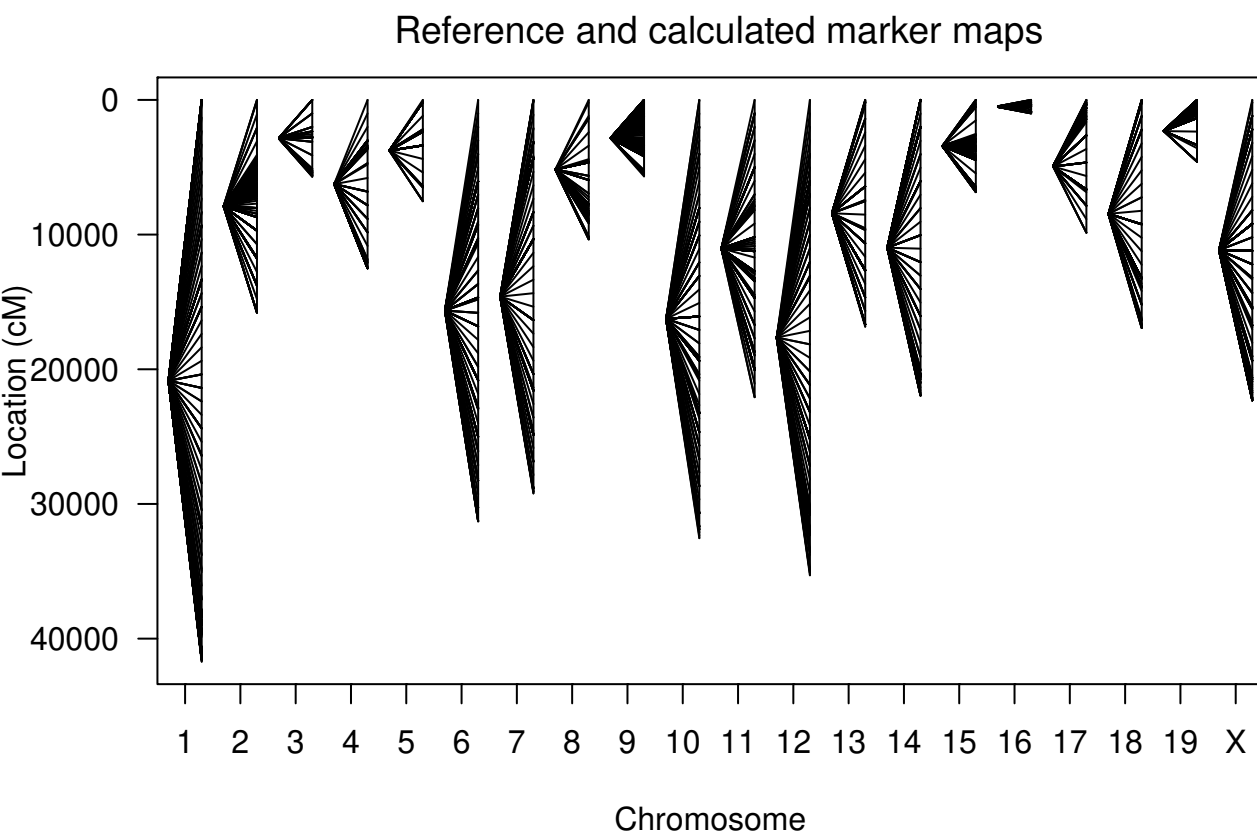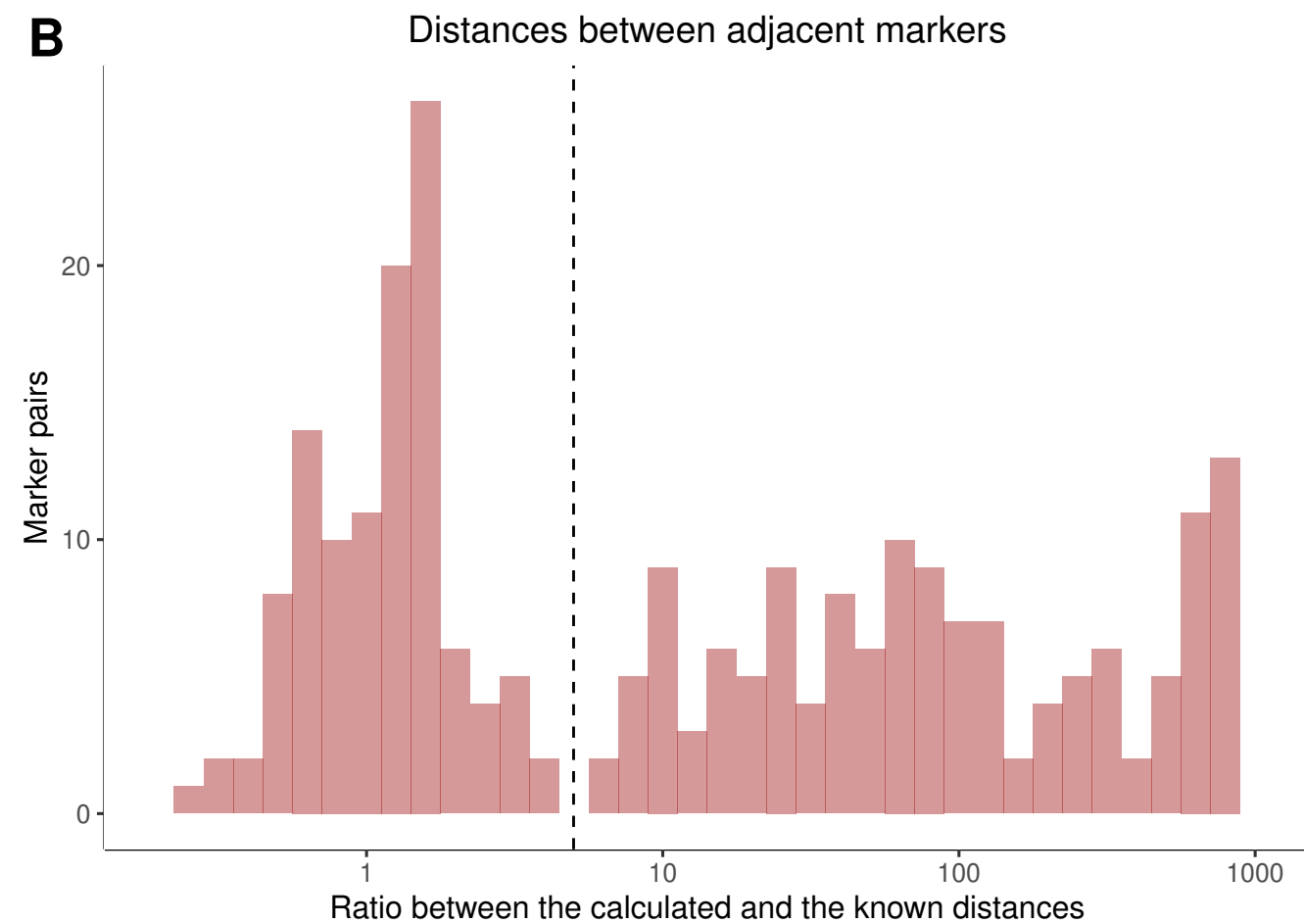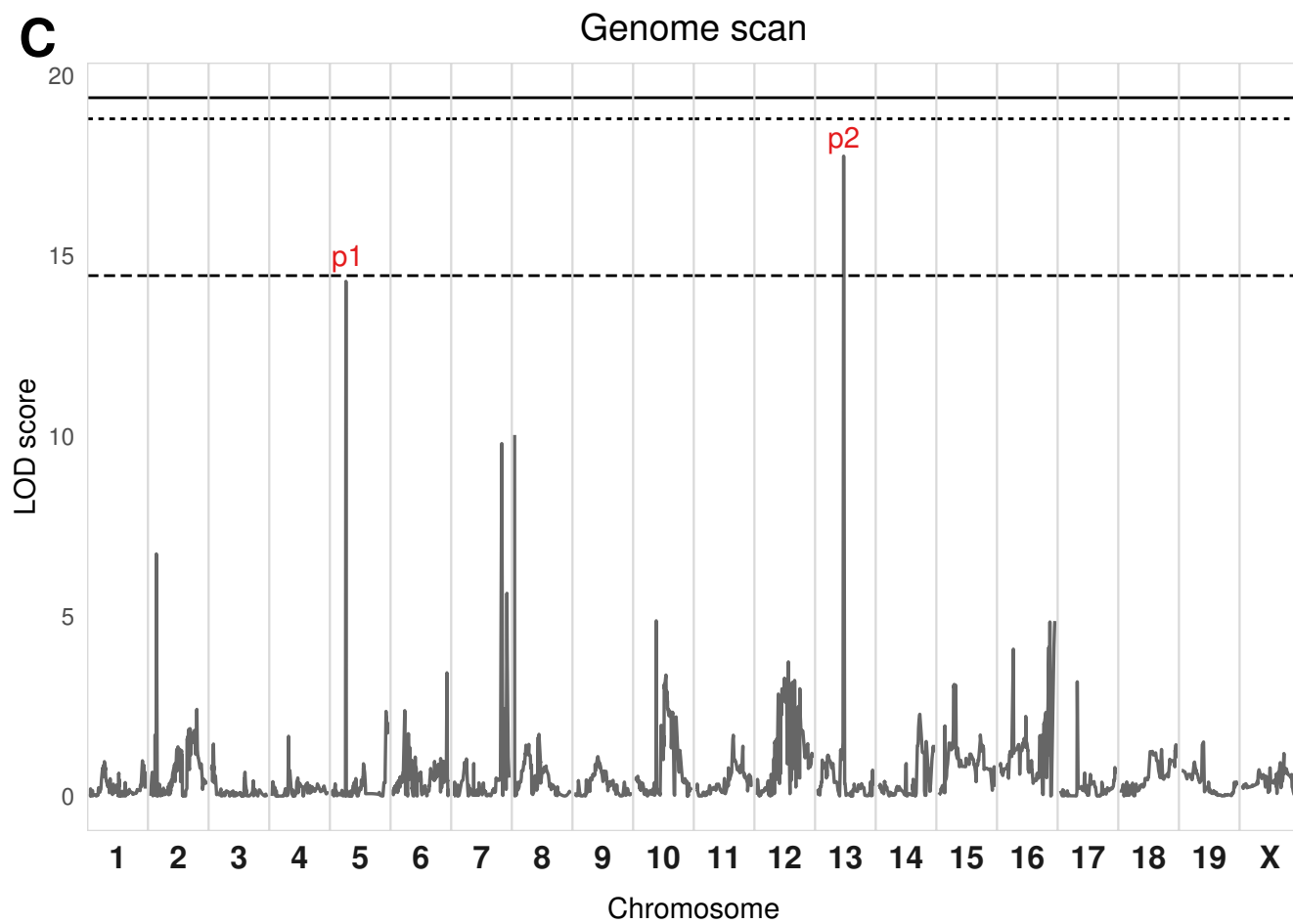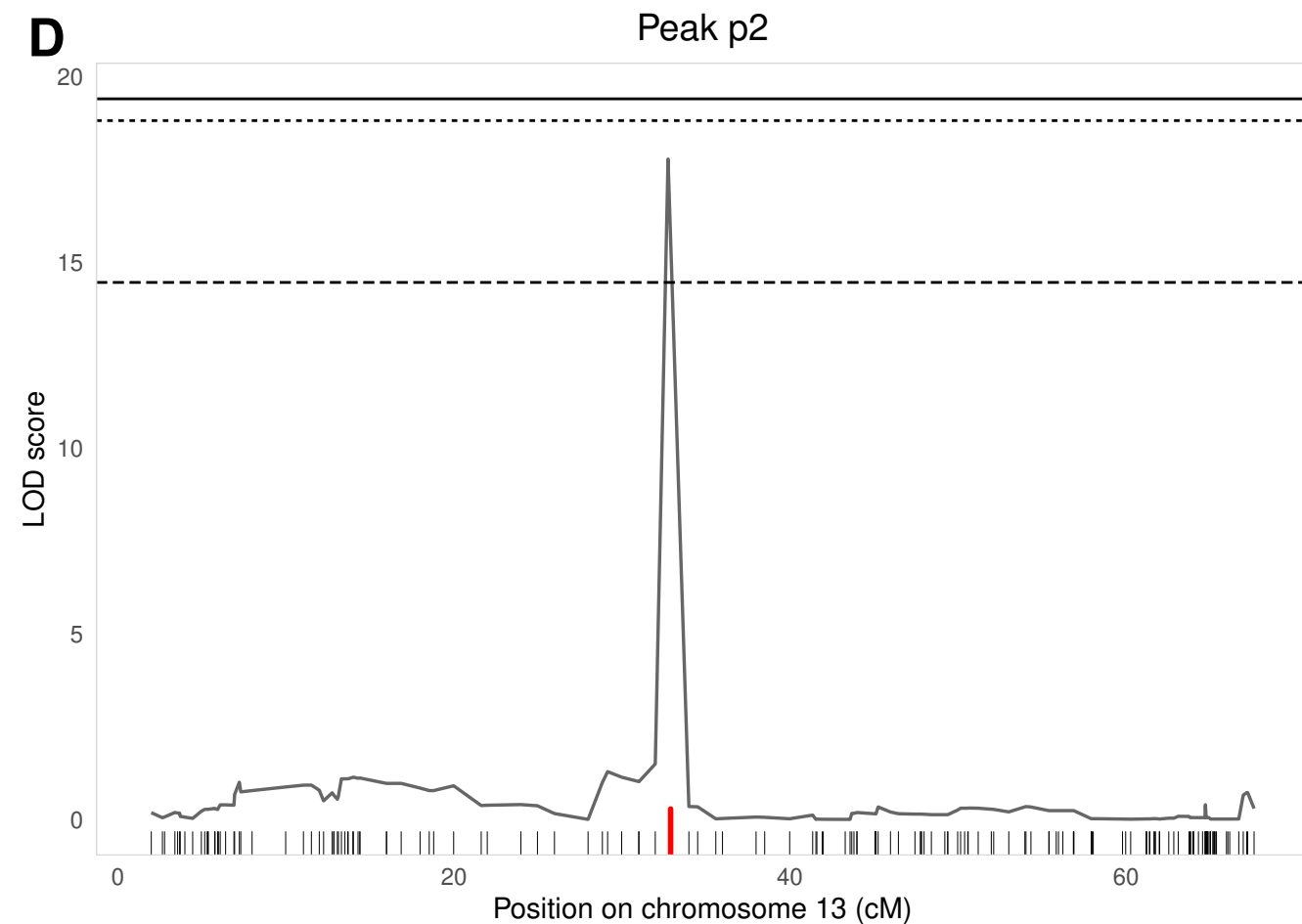394 **Supplementary Figure 4**

395 Analysis of the estimated genetic map in datasets 2 (A) and 3 (B) after curation of genotyping data by
396 stuart. Refer to Supplementary Figure 2 for comparison with original data. The estimated marker
397 maps are now consistent with the known marker maps with similar genome length despite local
398 contractions and expansions (the ratio between the calculated and the known length of the genome
399 is 1.00 for dataset 2 and 0.96 for dataset 3). The ratios between estimated and known distance
400 between adjacent markers in dataset 2 (C) and dataset 3 (D) are now normally distributed with a
401 mean=1.27 and a sd=0.77 for dataset 2 and a mean=1.24 and a sd=0.61 for dataset 3. The x-axis is in
402 logarithmic                                                                                      scale.

403

404 **Supplementary Figure 5**

405 Analysis of the LOD score curve after curation of genotyping data by stuart in datasets 2 (A) and 3
406 (B). Refer to Supplementary Figure 3 for comparison with original data. Significance thresholds are
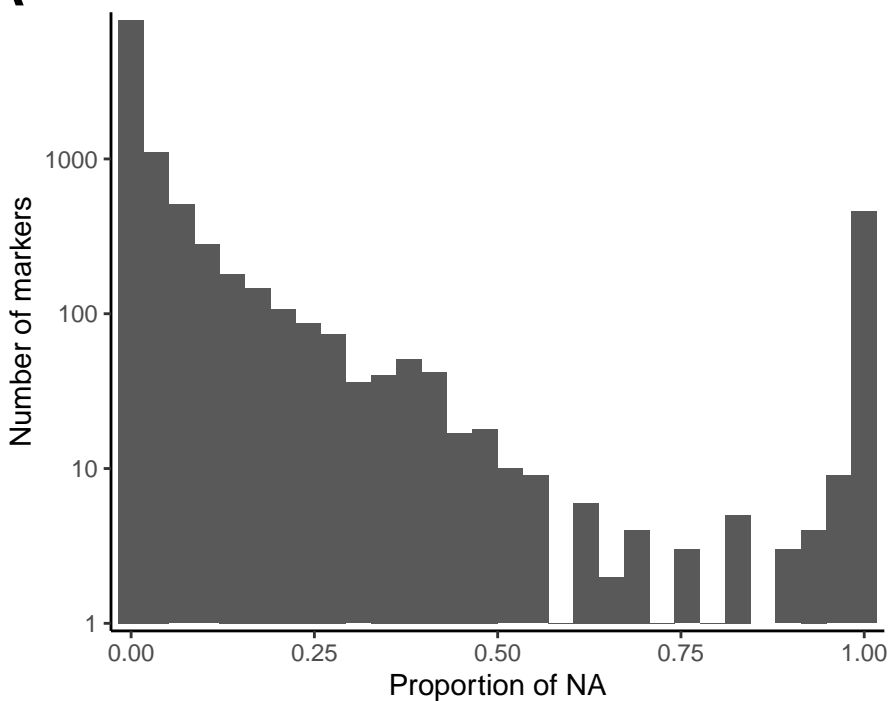
14

407     much lower than before curation. One peak in dataset 2 is significant at P<0.05 (plain line) and none

408     of the peaks observed before data curation (Supplementary Figure 3) were confirmed after curation

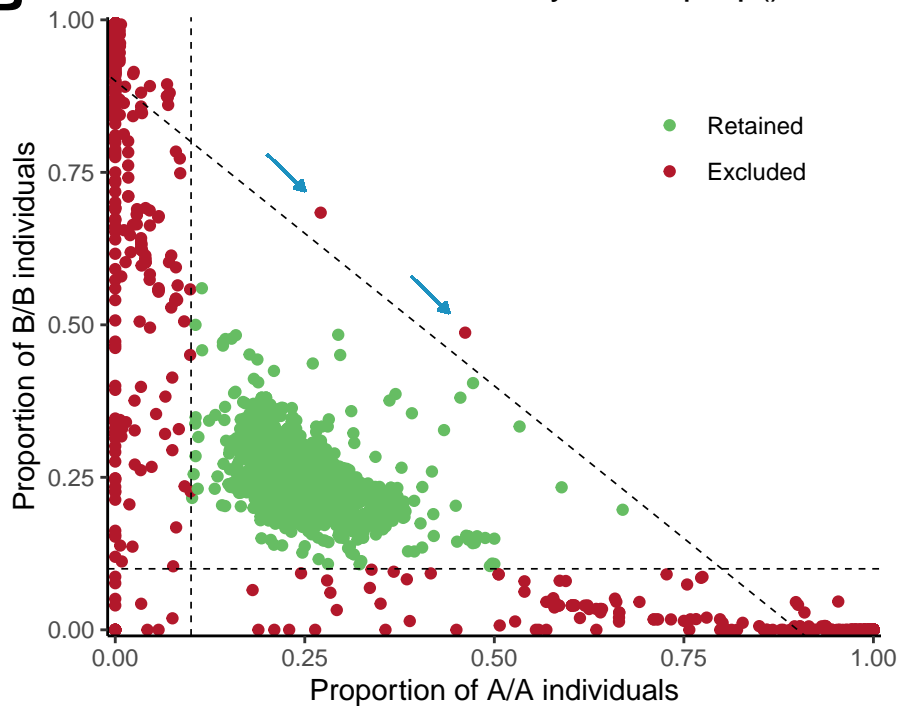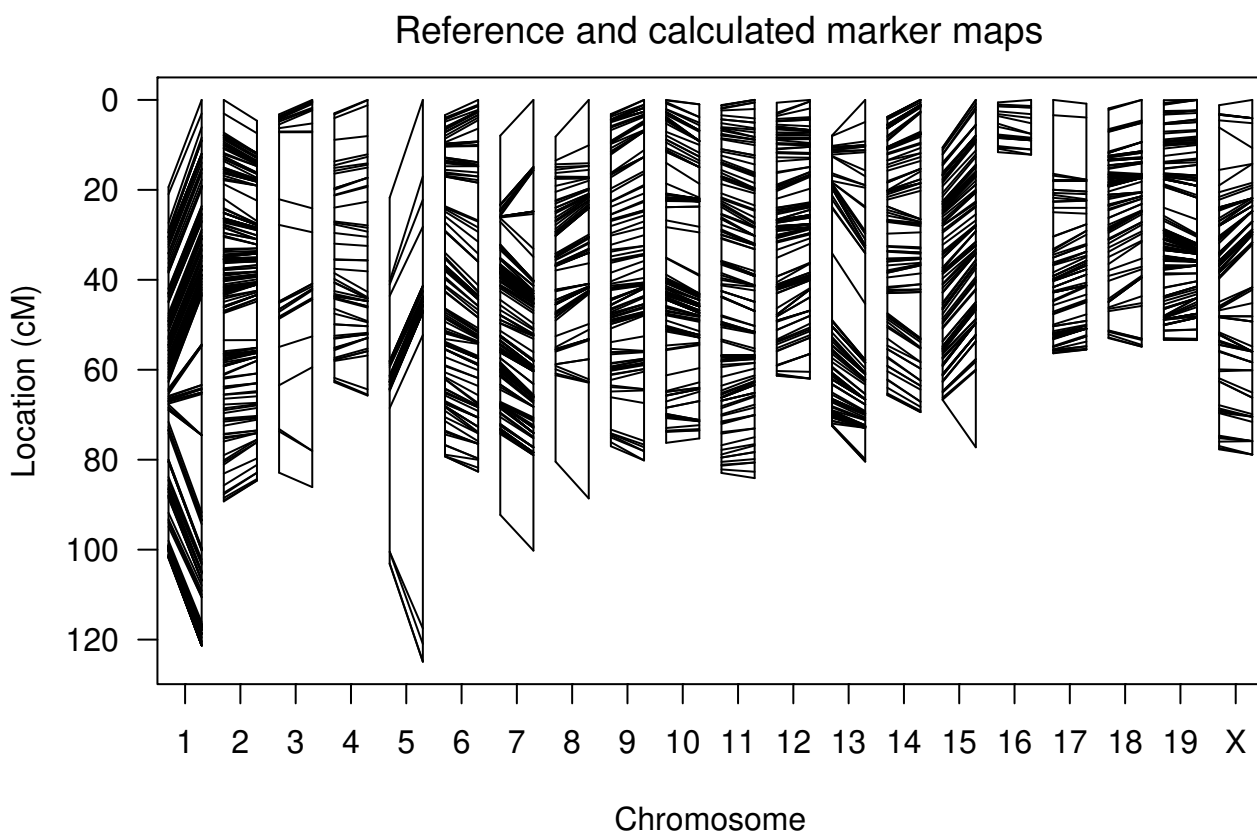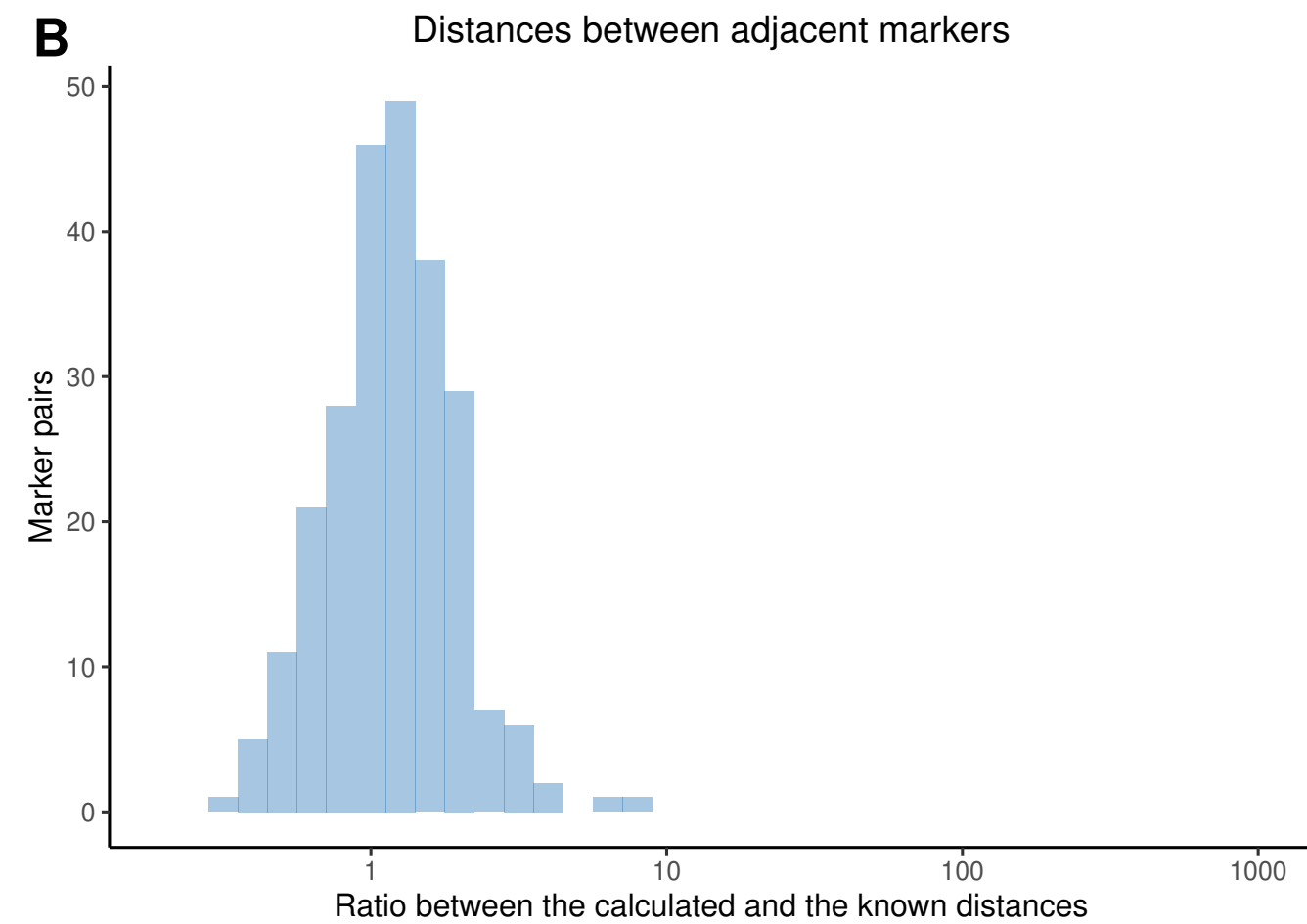409     with stuart. Dotted line: P=0.1. Dashed line: P=0.63.

410

**A** Reference and calculated marker maps

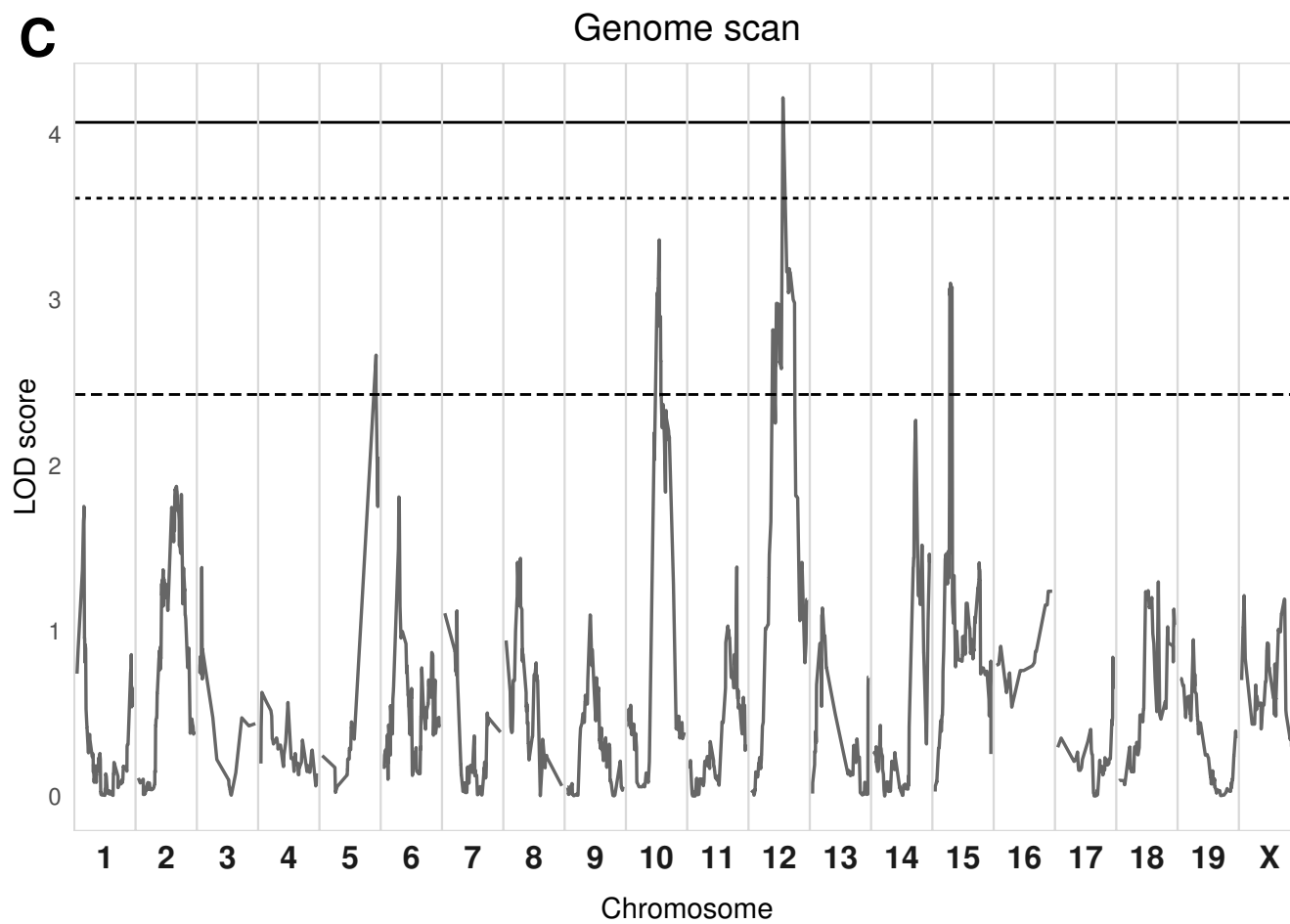**B** Distances between adjacent markers

**C** Genome scan

**D** Peak p2

**A** Proportion of missing genotypes

**B** Exclusion of markers by mark_prop()

**A** Reference and calculated marker maps

**B** Distances between adjacent markers

**C** Genome scan

**D** Chromosome 12 QTL