

Determinants of QTL mapping power in the realized Collaborative Cross

Gregory R. Keele^{*,†,‡}, Wesley L. Crouse^{*,†,‡}, Samir N. P. Kelada^{‡,§} and William Valdar^{‡,**,1}

*Authors contributed equally, [†]Curriculum in Bioinformatics and Computational Biology, [‡]Department of Genetics, [§]Marsico Lung Institute, ^{**}and Lineberger Comprehensive Cancer Center, University of North Carolina, Chapel Hill, North Carolina 27599
ORCID IDs: 0000-0002-1843-7900 (G.R.K.), 0000-0001-5745-4490 (W.L.C.), 0000-0003-2676-9232 (S.N.P.K.), 0000-0002-2419-0430 (W.V.)

ABSTRACT The Collaborative Cross (CC) is a mouse genetic reference population whose range of applications includes quantitative trait loci (QTL) mapping. The design of a CC QTL mapping study involves multiple decisions, including which and how many strains to use, and how many replicates per strain to phenotype, all viewed within the context of hypothesized QTL architecture. Until now, these decisions have been informed largely by early power analyses that were based on simulated, hypothetical CC genomes. Now that more than 50 CC strains are available and more than 70 CC genomes have been observed, it is possible to characterize power based on realized CC genomes. We report power analyses based on extensive simulations and examine several key determinants of power: 1) the number of strains and biological replicates, 2) the QTL effect size, and 3) the distribution of functionally distinct alleles among the founder strains at the QTL. We also provide general power estimates to aide in the design of future experiments. All analyses were conducted with our R package, SPARCC (Simulated Power Analysis in the Realized Collaborative Cross), developed for performing either large scale power analyses or those tailored to particular CC experiments.

KEYWORDS recombinant inbred lines, haplotype association, allelic series, multiparental population, MPP, quantitative trait, complex trait

Introduction

The Collaborative Cross (CC) is a multiparental population (MPP) recombinant inbred (RI) strain panel of laboratory mice derived from eight inbred founder strains (letter abbreviation in parentheses): A/J (A), C57BL/6J (B), 129S1/SvImJ (C), NOD/ShiLtJ (D), NZO/H1LtJ (E), CAST/EiJ (F), PWK/PhJ (G), and WSB/EiJ (H) (Threadgill *et al.* 2002; Churchill *et al.* 2004; Chesler *et al.* 2008; Threadgill and Churchill 2012). This set of founder strains represents three subspecies of the house mouse *Mus musculus* (Yang *et al.* 2011) and, in large part due to the inclusion of three wild-derived founders (F-H), imbues the CC panel with far greater genetic variation than previous RI panels derived solely from pairs of classical inbred strains. As an RI panel, the CC thus provides a diverse set of reproducible genomes and represents a powerful tool for genetic analysis

(Collaborative Cross Consortium 2012; Srivastava *et al.* 2017). Indeed, although the CC RI panel has only become available in the last six years (Welsh *et al.* 2012), it has already yielded new insights into human disease and basic mouse biology (Shusterman *et al.* 2013; Rogala *et al.* 2014; Rasmussen *et al.* 2014; Lorè *et al.* 2015; Gralinski *et al.* 2015; Venkatratnam *et al.* 2017; Orgel *et al.* 2018).

As originally envisaged, a key use of the CC was as a resource for QTL mapping (Threadgill *et al.* 2002; Churchill *et al.* 2004). In theory, its broad genetic diversity makes it ideal for this purpose, and its replicability permits the mapping of phenotypes such as drug-response that are otherwise hard to measure in outbreds (Mosedale *et al.* 2017). Its utility for QTL mapping in practice was also predicted by studies in the incipient CC lines, the pre-CC (Aylor *et al.* 2011; Durrant *et al.* 2011; Philip *et al.* 2011; Mathes *et al.* 2011; Kelada *et al.* 2012; Ferris *et al.* 2013; Ram *et al.* 2014; Rutledge *et al.* 2014; Kelada 2016; Donoghue *et al.* 2017; Phillippi *et al.* 2014)

Nonetheless, QTL mapping power depends in part on the number of strains available, and the number strains available in the CC is, and will remain, far less than the 1,000 proposed in Churchill *et al.* (2004): At the time of this work, mice were

Manuscript compiled: Friday 9th November, 2018

¹Corresponding author: 120 Mason Farm Rd., Genetic Medicine Bldg., Suite 5022, Campus Box 7264, University of North Carolina, Chapel Hill, NC 27599.

E-mail: william.valdar@unc.edu

available for 59 CC strains from the UNC Systems Genetics Core, with a subset from these 59 and an additional 11 expected to be offered through the Jackson Laboratory (JAX), a total of 70 CC strains potentially.

A reduction in strain numbers as a function of allelic incompatibilities between subspecies (Shorter *et al.* 2017) was expected, and winnowed the number of resulting CC strains down to 50-70. Although smaller than originally intended, this population size reflects the biological and financial realities of maintaining a sustainable mammalian genome reference population. [Whereas cost grows proportional to the number of strains, demand does not, and a much larger number of strains would threaten the economic viability of the operation (F. Pardo-Manuel de Vilena, *pers. comm.*.)] Nonetheless, subsets of the available CC strains have already been used to map QTL, as evidenced by a growing list of studies (Vered *et al.* 2014; Mosedale *et al.* 2017; Graham *et al.* 2017). Beyond these successes, however, it is unclear how much the reduction has affected the ability to map QTL in the CC in general.

The initially proposed figure of 1,000 CC strains in Churchill *et al.* (2004) was more formally justified in Valdar *et al.* (2006a) as being necessary to provide enough power both to map single QTL and for robust, genome-wide detection of epistasis. That estimate was based on simulations involving larger numbers (500-1,000) of hypothetical CC genomes. Those simulations, performed before any CC strains existed and with the goal of guiding the CC's design, had a broad scope, exploring the effect of varying strain numbers, alternative mapping approaches [association of single nucleotide polymorphisms (SNPs) vs association of inferred haplotypes], and alternative breeding strategies. As such, the power estimates that were reported do not reflect the number of CC strains now available, nor their actual, realized founder mosaic genomes. An updated, more focused power analysis that both exploits and works within the constraints of the realized genomes is therefore timely.

Power analyses have been performed previously for a number of RI panels. For biparental RIs, they have been performed analytically in plants (*e.g.*, Kaepler 1997), animals [*e.g.*, the BXD lines in mice (Belknap *et al.* 1996; Peirce *et al.* 2004)], and in general (Cowen 1988; Soller and Beckmann 1990; Knapp and Bridges 1990), as well as through simulation (Falke and Frisch 2011; Takuno *et al.* 2012). For MPP RIs, they have most often been reported as those resources are introduced to the community. This includes, in plants: *Arabidopsis* (Kover *et al.* 2009; Klaseen *et al.* 2012), nested association mapping (NAM) populations (Li *et al.* 2011) in maize (Yu *et al.* 2008) and sorghum (Bouchet *et al.* 2017), and multigenerational advanced intercross (MAGIC) populations of rice (Yamamoto *et al.* 2014) and maize (Dell'Acqua *et al.* 2015). In animals, other than aforementioned prospective study of Valdar *et al.* (2006a): Noble *et al.* (2017) assessed mapping power of SNP association while introducing a 507-strain nematode resource, the *Caenorhabditis elegans* Multiparental Experimental Evolution (CeMEE) panel; and King *et al.* (2012) estimated haplotype-based association power while introducing the *Drosophila* Synthetic Population Resource (DSPR), a fly panel with more than 1,600 lines. In a follow-up DSPR power analysis, King and Long (2017) compared the DSPR with the related *Drosophila* Genetic Reference Panel (DGRP) (Mackay *et al.* 2012). They illustrated how QTL effect size differs between a population whose allele frequencies are balanced (DSPR) vs one whose allele frequencies are less balanced (DGRP) and explored implications for cross-population validation; they also compared

mapping power for bi-allelic QTL, based on single SNPs, and multi-allelic QTL constructed from actual adjacent SNPs within genes.

Here we examine related topics on QTL mapping power in the realized CC, including: 1) how power is affected by the number of strains and replicates; 2) how it is affected by the number of functional alleles and their distributions among the founders; and 3) how the QTL effect size is specific to a particular population or sample and how that influences a power estimate's interpretation.

To allow researchers to repeat our analyses, but tailored to their own specific requirements or with updated CC genome lists, we provide an R package SPARCC (Simulated Power Analysis of the Realized Collaborative Cross), a tool that evaluates the power to map QTL by performing efficient haplotype regression-based association analysis of simulated QTL using the currently available CC genomes. SPARCC is highly flexible, allowing QTL to be simulated with any possible allele-to-founder pattern and scaled with respect to different reference populations. As a reusable resource, researchers could estimate power calculations based on the CC strains available to them and potentially incorporate prior knowledge about the genetic architecture of the likely QTL or the phenotype as whole.

Methods

Our power calculations are based on three main processes:

1. Simulation of CC data, including selection of CC strains from a fixed set of realized CC genomes, and QTL location, and simulation of phenotypes.
2. QTL mapping, including determination of significance thresholds.
3. Evaluation of QTL detection accuracy, power and false positive rate.

These are described in detail below, after a description of the genomic data that serves as the basis for the simulations.

Data on realized CC genomes

CC strains. Genome data was obtained for a set of 72 CC strains (listed in **Appendix D**) available at the time of writing from <http://csbio.unc.edu/CCstatus/index.py?run=FounderProbs>. Genome data was in the form of founder haplotype mosaics (see below) for each strain, this based on genotype data from the MegaMUGA genotyping platform (Morgan *et al.* 2016) applied to composites of multiple mice per strain. Although some of the 72 strains will likely become extinct (Darla Miller *pers. comm.*), it is also possible that more may be added. At the time of writing, however, these were all genomes that had been observed.

The availability of mice from the 72 strains is as follows. The UNC Systems Genetics Core is currently planning to maintain and distribute 59 strains (also listed in **Appendix D**). The Jackson Laboratory will eventually maintain some subset of these same 59 strains and potentially an additional 11 strains not maintained at UNC. For five of these currently available 59 strains, founder haplotype mosaics were not available on the website at the time of this study, and thus not included in these simulations. (Note: The mosaics became available at the time of writing, and could be included in future simulations; see **Discussion**.) Two of the strains included in both these simulations and the 59 available strains were derived from the same breeding funnel (CC051

and CC059); because of their close relatedness, which was not explicitly modeled in the simulations, we emphasize that there are 58 independent strains that are currently available, and include an indicator in figures denoting this number as a currently realistic maximum for strain number in CC studies.

Reduced dataset of haplotype mosaics. The genomes of the CC, as with other MPPs, can be represented by inferred mosaics of the original founder haplotypes (Mott *et al.* 2000). Founder haplotype mosaics were inferred previously by the UNC Systems Genetics Core (<http://csbio.unc.edu/CCstatus/index.py?run=FounderProbs>) using the hidden Markov model (HMM) of Fu *et al.* (2012) applied to genotype calls from MegaMUGA, a genotyping platform providing data on 77,800 SNP markers (Morgan *et al.* 2016). The HMM inference provides a vector of 36 diplo-type probabilities for each CC strain for each of 77,551 loci (each defined as the interval between adjacent, usable SNPs) across the genome. Rather than using all of the available data for our simulations, we used a reduced version: since adjacent loci often have almost identical descent, mapping using all loci is both computationally expensive and—at least for the purposes of the power analysis—largely redundant. Thus, prior to analysis the original dataset was reduced by averaging adjacent genomic intervals whose diplo-type probabilities were highly similar. Specifically, adjacent genomic intervals were averaged if the maximum L2 norm between the probability vectors of all individuals is less than 10% of the maximum possible L2 norm ($\sqrt{2}$); this reduced the file storage from 610 MB to 288 MB, and the genome from 77,551 to 17,900 intervals (76.9% reduction in positions to be evaluated in a scan).

Phenotype simulation

Phenotypes for CC strains were simulated based on effects from a single QTL, plus effects of polygenic background (“strain effects”), and noise. Within our simulation framework, we specified: 1) the QTL location, which randomly was sampled from the genome; 2) the sample size in terms of both strains and replicates; 3) how the eight possible haplotypes at that location are grouped into eight or fewer functional alleles (the “allelic series”; see below); and 4) how those alleles, along with strain information, are used to generate phenotype values (see below).

Underlying phenotype model. Simulated phenotypes were generated according to the following linear mixed model. For given QTL with $m \leq 8$ functional alleles, phenotype values $y = \{y_i\}_{i=1}^N$ for N individuals in $n \leq N$ strains were generated so that

$$y = \mathbf{1}\mu + \underbrace{\mathbf{Z}\mathbf{X}\boldsymbol{\beta}}_{\text{QTL effect}} + \underbrace{\mathbf{Z}\mathbf{u}}_{\text{Strain effect}} + \underbrace{\boldsymbol{\varepsilon}}_{\text{Noise}}, \quad (1)$$

where $\mathbf{1}$ is an N -vector of 1’s, μ is an intercept, \mathbf{Z} is an $N \times n$ incidence matrix mapping individuals to strains, \mathbf{X} is an $n \times m$ allele dosage matrix mapping strains to their estimated dosage of each of the m alleles, $\boldsymbol{\beta}$ is an m -vector of allele effects, \mathbf{u} is an n -vector of strain effects (representing polygenic background variation), and $\boldsymbol{\varepsilon}$ is an N -vector of unstructured, residual error. The parameter vectors $\boldsymbol{\beta}$, \mathbf{u} and $\boldsymbol{\varepsilon}$ were each generated as being equivalent to independent normal variates rescaled to have specific variances: the strain effects \mathbf{u} and residual $\boldsymbol{\varepsilon}$ were rescaled to have population (rather than sample) variances h_{strain}^2 and σ^2 respectively; the allele effects $\boldsymbol{\beta}$ were rescaled so that the QTL contributes a variance h_{QTL}^2 , with this latter rescaling performed in one of three distinct ways (described later).

The relative contributions of the QTL, polygenic background, and noise were thus controlled through three parameters: the QTL effect size, h_{QTL}^2 , the strain effect size, h_{strain}^2 , and the residual variance σ^2 . By convention, these were specified as fractions summing to exactly 1.

The allele dosage matrix \mathbf{X} was generated by collapsing functionally equivalent haplotypes according to a specified allelic series. Let \mathbf{D} be an $n \times 36$ incidence matrix describing the haplo-type pair (diplotype) state of of each CC strain at the designated QTL, with columns corresponding to AA, ..., HH, AB, ..., GH, such that, for example, $\{D\}_{3,1} = 1$ implies CC strain 3 has diplotype AA. Then

$$\mathbf{X} = \mathbf{D}\mathbf{A}\mathbf{M}, \quad (2)$$

where \mathbf{A} is an 36×8 additive model matrix that maps diplotype state to haplotype dosage (*e.g.*, diplotype AA equals 2 doses of A), and \mathbf{M} is an $8 \times m$ “merge matrix” [after Yalcin *et al.* (2005)] that encodes the allelic series, mapping the 8 haplotypes to m alleles, such that if haplotypes A and B were both in the functional group “allele 1”, then diplotype AB in \mathbf{D} would correspond to 2 doses of allele 1 in \mathbf{X} (see examples in **Appendix E**).

QTL allelic series. The specification of an allelic series, rather than assuming all haplotype effects are distinct, acknowledges that for many QTL we would expect the same functional allele to be carried by multiple founder haplotypes. For our main set of simulations, the allelic series was randomly sampled from all possible configurations (examples in **Figure 1**); in a smaller, more focused investigation of the effects of allele frequency imbalance, we sampled from all possible configurations of bi-alleles.

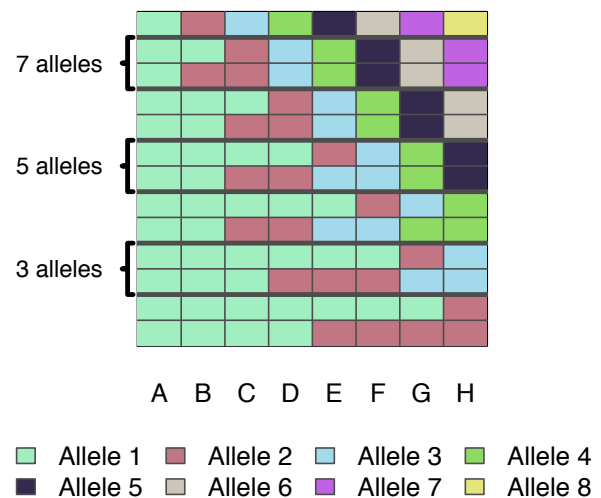


Figure 1 Example allelic series with differing numbers of functional alleles. Each row is an allelic series, each column of the grid is a CC founder, and colors correspond to functional allele. Two examples of allelic series are provided for each number of functional alleles: a balanced series and an imbalanced series. The entire space of allelic series are not shown here; however, the full space of series with two alleles is shown in **Figure 7A**.

Alternative definitions of QTL effect size: B and DAMB. The QTL effect size (h_{QTL}^2) is a critical determinant of mapping power; yet its precise definition and its corresponding interpretation often varies between studies and according to what

question is being asked. We used two alternative definitions, “B” and “DAMB”, described below. These alternatives acknowledge that the proportion of variance explained by a particular QTL, and thus the power to detect that QTL, is not determined solely by h_{QTL}^2 , but rather depends on several additional factors, namely: the variance of the finite sample of allele effects β ; the allelic series configuration \mathbf{M} ; and the particular set of CC strains and their locus diplotypes \mathbf{D} .

Definition B scales the allele effects so that $h_{\text{QTL}}^2 = V(2\beta)$, where $V()$ denotes the population variance (rather than the sample variance). The QTL effect size is interpretable as the variance that would be explained by the QTL in a theoretical population that is balanced with respect to the functional alleles. As such, the proportion of variance explained by the QTL in the mapping population will deviate from h_{QTL}^2 due to imbalance in both \mathbf{M} and \mathbf{D} . Conversely, for a given h_{QTL}^2 , the allelic values at a QTL will be constant across populations. (Note: the 2 multiplier ensures proper scaling since \mathbf{X} from Eq 2 includes dosages of founder haplotypes at the QTL, ranging from 0 to 2.)

Definition DAMB scales the QTL effect so that $h_{\text{QTL}}^2 = V(\mathbf{DAM}\beta)$. The QTL effect size is exactly the variance explained by the QTL in the mapping population, essentially the R^2 . As such, it depends on both \mathbf{M} and \mathbf{D} . Correspondingly, for a given h_{QTL}^2 , the allelic values will adjust depending on which population they are in. [In Appendix A, for completeness, we also describe a further, intermediate option, Definition MB, where $h_{\text{QTL}}^2 = V(2\mathbf{M}\beta)$, corresponding to balanced founder contributions.]

The earlier power study of Valdar *et al.* (2006a), which considered only bi-allelic QTL, defined effect size in a manner comparable to Definition B.

Averaging over strains and causal loci. The previous subsections described simulation of a single phenotype conditional on a set of strains and a causal genomic locus. For each of S simulations, $s = 1, \dots, S$, we averaged over these variables by uniformly sampling 1) the set of strains included in the experiment (for a specified number of strains), 2) the causal locus underlying the QTL, and 3) the allelic series (for a specified number of functional alleles). This was intended to produce power estimates that take into account many sources of uncertainty and are thus broadly applicable.

QTL detection and power estimation

QTL mapping model. QTL mapping of the simulated data was performed using a variant of Haley-Knott (HK) regression (Haley and Knott 1992; Martínez and Curnow 1992) that is commonly used in MPP studies (Mott *et al.* 2000; Liu *et al.* 2010; Fu *et al.* 2012; Gatti *et al.* 2014; Zheng *et al.* 2015) whereby association is tested between the phenotype and the local haplotype state, the latter having been inferred probabilistically from genotype (or sequence data) and represented as a set of diplotype probabilities or, in the case of an additive model, a set of haplotype dosages then used as predictors in a linear regression. Specifically, we used HK regression on the strain means (Valdar *et al.* 2006a; Zou *et al.* 2006) via the linear model

$$\bar{\mathbf{y}}^{(s)} = \mathbf{1}\mu + \mathbf{P}\mathbf{A}\beta + \epsilon, \quad (3)$$

where $\bar{\mathbf{y}}^{(s)}$ is the s^{th} simulated n -vector of strain means, \mathbf{P} is an $n \times 36$ matrix of inferred diplotype probabilities for the sampled CC genomes at the QTL [i.e., $\mathbf{P} = p(\mathbf{D}|\text{genotype data})$]; see Zhang *et al.* (2014)], and ϵ is the n -vector of residual error on

the means, distributed as $\epsilon \sim N(\mathbf{0}, \mathbf{I}(h_{\text{strain}}^2 + \sigma^2/r))$. The above implies an eight-allele model (cf Equation 1 with $\mathbf{M} = \mathbf{I}$). Although this could lead to reduced power when there are fewer functional alleles, particularly at loci in which the functional alleles are not well represented, it is most common in practice, in accordance with the fact that the allelic series of an unmapped QTL would typically be unknown in advance [e.g., Mott *et al.* (2000); Valdar *et al.* (2006a,b); Svenson *et al.* (2012); Gatti *et al.* (2014)]. The fit of Eq 3 was compared with that of an intercept-only null model via an F-test, and produced a p-value, reported as its negative base 10 logarithm, the logP. This procedure was performed for all loci across the genome, resulting in a genome scan for $\mathbf{y}^{(s)}$.

Genome-wide significance thresholds and QTL detection.

Genome-wide significance thresholds were determined empirically by permutation. The CC panel is a balanced population with respect to founder genomic contributions and, by design, has minimal population structure. These features support the assumption of exchangeability among strain genomes: that under a null model in which the genetic contribution to the phenotype is entirely driven by infinitesimal (polygenic) effects, all permutations of the strain labels (or equivalently, of the strain means vector $\mathbf{y}^{(s)}$) are equally likely to produce a given configuration of $\mathbf{y}^{(s)}$. Permutation of the strain means, $\mathbf{y}^{(s)}$, was therefore used to find the logP critical value controlling genome-wide type I error rate (GWER) (Doerge and Churchill 1996). Briefly, we sampled p permutations and perform genome scans for each; this was done efficiently using a standard matrix decomposition approach (Appendix B). The maximum logPs per genome scan and simulation s were then recorded, and these are fitted to a generalized extreme value distribution (GEV) (Dudbridge and Koeleman 2004; Valdar *et al.* 2006a) using R package *evir* (Pfaff and McNeil 2018). The upper $\alpha = 0.05$ quantile of this fitted GEV was then taken as the α -level significance threshold, $T_{\alpha}^{(s)}$. If the maximum observed logP for $\mathbf{y}^{(s)}$ in the region of the simulated QTL exceeded $T_{\alpha}^{(s)}$, then the corresponding locus was considered to be a (positively) detected QTL (see immediately below).

Performance evaluation. For a given simulation, we declared a true positive if the detected QTL was within ± 5 Mb of the true (simulated) QTL. The 5 Mb window size was used to approximate a QTL support interval, which is partly a function of linkage disequilibrium (LD) in the CC. (LD has been characterized in the CC previously but not summarized with a single point estimate (Collaborative Cross Consortium 2012); our choice of 5 Mb is therefore an approximation, but we find that it only marginally increased mapping power relative to using smaller window widths.) A false positive was declared if a QTL was detected on a chromosome other than that harboring the QTL. Simulations in which a QTL was detected on the correct chromosome but outside the 5 Mb window were disregarded; although potentially wasteful of data, this measure avoided the arbitrariness of formulating rules for edge cases in which it was ambiguous whether the simulated signal was detected or not. Power for a given simulation setting was then defined as the proportion of true positives among all simulations at that setting, and the false positive rate was defined as the proportion of false positives.

As a measurement of mapping resolution, for true positive detection, we recorded the mean and the 95% quantile of the

genomic distance from the true QTL. Given our criterion for calling true positives, the maximum distance was necessarily 5 Mb, and experimental settings that correspond to low power would be expected to have fewer data points, yielding estimates that are unstable. In order to obtain more stable estimates, we used a regularization procedure, estimating the mean distance and 95% quantiles as weighted averages of the observed values and prior pseudo-observations. Specifically, for an arbitrarily small but detected true positive QTL, it is reasonable to expect the peak signal to be distributed uniformly within the ± 5 Mb window. This implies a mean location error of 2.5 Mb and a 95% quantile of 4.75 Mb. Thus, when calculating the regularized mean location error we assumed 10 prior pseudo-observations of 2.5 Mb, and when calculating the regularized 95% quantile we assume 10 prior pseudo-observations of 4.75 Mb. This number of pseudo-observations represents 1% of the maximum number of possible data points.

Overview of the simulations

Simulation settings. Simulations for all combinations of the following parameter settings:

- Number of strains: [(10-70 by 5), 72]
- QTL effect size (%): [1, (5-95 by 5)]
- Number of functional alleles: [2, 3, 8]

The number of observations per strain were fixed at $r = 1$ and the background strain effect size was fixed at $h^2_{\text{strain}} = 0\%$ with the understanding that results from these simulations provide information on other numbers of replicates and strain effect sizes implicitly. Specifically, a simulated mapping experiment on strain means that assumes r replicates, strain effect h^2_{strain} , and QTL effect size h^2_{QTL} is equivalent to a single-observation mapping experiment with no strain effect and QTL effect size \bar{h}^2_{QTL} , where

$$\bar{h}^2_{\text{QTL}} = \frac{h^2_{\text{QTL}}}{h^2_{\text{QTL}} + h^2_{\text{strain}} + \sigma^2/r} \quad (4)$$

[Valdar *et al.* (2006a), after Soller and Beckmann (1990); Knapp and Bridges (1990); Belknap (1998)]. For example, a mapping experiment on strain means with QTL effect size $h^2_{\text{QTL}} = 0.3$, $h^2_{\text{strain}} = 0.4$, $\sigma^2 = 0.3$, and $r = 10$, is equivalent to our simulation of a single-observation with no strain effect but QTL effect size $\bar{h}^2_{\text{QTL}} \simeq 0.41$ (**Appendix C**).

We conducted 1,000 simulations per setting. CC strains and the position of the QTL were sampled for each simulation, providing estimates of power that are effectively averaged over the CC population. We ran these settings for QTL effect sizes specified with respect to the observed mapping population (Definition DAMB) and a theoretical population that is balanced in terms of the functional alleles (Definition B). Confidence intervals for power were calculated based on Jeffreys interval (Brown *et al.* 2001) for a binomial proportion. A description of the computing environment and run-times are provided in **Appendix C**.

Availability of data and software

R package. All analyses were conducted in the statistical programming language R (R Core Team 2018). SPARCC is available as an R package on GitHub at <https://github.com/gkeele/sparcc>. Specific arguments that control the phenotype simulations, the

strains used, genomic position of simulated QTL, and allelic series, are listed in **Appendix A**.

Also included with the SPARCC R package are several results datasets. These include data tables of power summaries from our simulations, as well as table summaries from simulations of a bi-allelic QTL that is balanced in the founders, maximally unbalanced in the founders, and the distance between detected and simulated QTL. Supplemental files for generating the simulated data, analyses, and figures are available at FigShare.

CC strains. The 72 CC strains with available data that were included in the simulations are described in **Appendix D**. Founder diplotype probabilities for each CC strain are available on the CC resource website (<http://csbio.unc.edu/CCstatus/index.py?run=FounderProbs>). We used probabilities corresponding to build 37 of the mouse genome, though build 38 is also available at the same website.

We store the founder haplotype data in a directory structure that SPARCC is designed to use, and was initially established by the HAPPY software package (Mott *et al.* 2000). The reduced data are available on GitHub at https://github.com/gkeele/sparcc_cache.

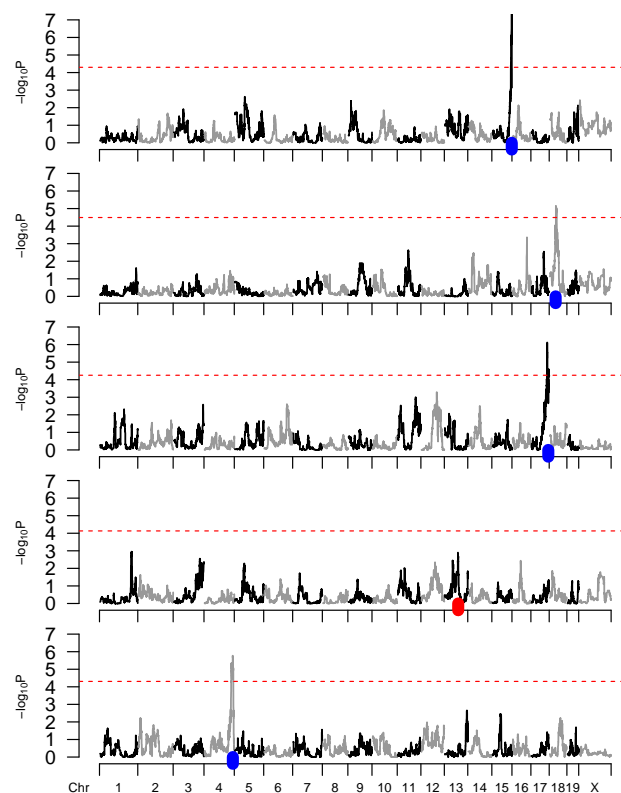


Figure 2 Simulated CC data and resulting genome scans. Five simulated genome scans are generated by the code provided in a simple example using our package SPARCC. Red dashed lines represent 95% significance thresholds based on 100 permutation scans. A blue tick represents the simulated position for a QTL that was successfully detected, whereas a red tick marks a QTL that was missed. These simulations were based on a specified set of 65 CC strains, five replicates of each strain, two functional alleles, 10% QTL effect size, and no background strain effect. The QTL is not mapped in the fourth simulation, ranked top to bottom, resulting in a power of 80%. Actual power calculations are based on a greater number of simulations.

Results

Power simulations were performed for varying numbers of strains, replicates and functional alleles, and for a ladder of QTL effect sizes. QTL effect size was defined in two ways: as the variance explained in a hypothetical populations that is balanced with respect to the alleles (Definition B; see **Methods**), or as the variance explained in the realized population (Definition DAMB). In this section we focus on results using the first of these, Definition B, owing to its more consistent theoretical interpretation. Under that definition, plots of power against numbers of strains are shown in Figure 3, and power across a representative selection of conditions is shown in **Table 1**. For comparison, these numbers are also provided for simulations under Definition DAMB in Table S1. Throughout these simulations the false positive rate was controlled at the target 0.05 level (**Figure S2**).

Large effect QTL usually detected by 50 or more strains.

Studies without replicates and with large numbers of strains (>50) were estimated to be well-powered to detect large effect QTL (>40%) (**Figure 3 [top]**); detecting smaller effect QTL, however, requires many replicates. For example, using 50 strains, the power to detect a 20% effect-size QTL with a single replicate is near-zero; with 5 replicates it approaches 80%. Detecting QTL with effect sizes $\leq 10\%$ was challenging. For example, achieving 80% power to detect an effect size of 10% when all 72 CC strains were used required more than 5 replicates per strain (**Figure 3 [middle right]**). Detecting even smaller QTL would require higher numbers of replicates. Additionally, the background strain effect harshly reduced QTL mapping power of small effect QTL (**Figure 3 [bottom]**).

Additional strains improve power more than additional replicates.

We investigated the relationship between power and the total number of mice, evaluating whether power gains were greater with additional CC strains or additional replicate observations. Power was interpolated over a grid of values for number of replicates and total number of mice from simulations based on a single observation per strain (Figure 5). This showed that additional CC strains improved mapping power more than additional replicates; this is indicated by higher power values for lower numbers of replicates while holding number of mice constant (see Figure 5, bordered vertical section at 250 mice).

Location error of detected QTL.

To obtain an approximation of mapping resolution, for all true positive detections we recorded the location error, or the genomic distance between simulated and detected QTL. The mean and the 95% quantile of the location error are reported as stabilized estimates for different numbers of strains and QTL effect sizes, but averaged over all other conditions, in Figure 4. (The stabilization procedure is described in **Methods**; raw, unstabilized estimates provided Figure S3.) The location error statistics require careful interpretation: for a detection to be classed as a true positive it had to be within 5Mb of the simulated QTL; therefore, location error was artificially capped at 5Mb. Mediocre performance thus corresponds to when that location seems uniformly (and therefore arbitrarily) distributed over the $\pm 5\text{Mb}$ interval, that is, having a mean of 2.5Mb and a 95% quantile of 4.8Mb.

Location error was improved (reduced) by increasing the number of strains, increasing the QTL effect size, or both. In particular, as with power, location error was improved by increasing the number of strains even when while holding the total number of mice constant (**Figure S4**), consistent with mapping resolution being improved by an increased number of recombination events in the QTL region. Distributions of raw location error, stratified by levels of the number of strains, the number of functional alleles, and the QTL effect size can be found in Figure S5.

Allele frequency imbalance reduces power

For a fixed set of QTL allele effects, it is expected that power will always be greatest when allele frequencies are balanced. Accordingly, when QTL effect size was defined in terms of the variance that would be explained in a theoretical population with balanced allele frequencies (Definition B), deviations from balance in the mapping population inevitably reduce power (Figure 6A). This reduction in power under Definition B is most evident for bi-allelic QTL (pink), in which the potential imbalance in allelic series is most extreme, namely when a single founder carries one functional allele and the other seven possess the alternative allele (7v1).

Conversely, when the QTL effect size is defined in terms of variance explained in the mapping population (Definition DAMB, which is similar to an R^2 measure), power remains constant across different allelic series and degrees of balance. Although note that this definition carries with it the (possibly unrealistic) implication that allele effects vary depending what population they are in.

When averaged over many allelic series, QTL mapping power based on Definition B is reduced relative to Definition DAMB, with the greatest reduction occurring for bi-allelic QTL (**Figure 6 B**). Though this modest reduction in power may seem to suggest that simulating with respect to a balanced population (Definition B) versus the mapping population (Definition DAMB) is unimportant in terms of designing a robust mapping experiment in the CC, we reiterate the value of using Definition B. Specifically, simulating with respect to Definition DAMB is overly optimistic regarding mapping power for QTL with imbalanced allelic series.

We performed additional simulations to evaluate bi-allelic QTL in more detail, these being more prone to drastic imbalance under Definition B. All 127 possible bi-allelic series are visualized as a grid in Figure 7A, ordered from balance and high power to imbalance and low power. The corresponding power estimates are shown in Figure 7B. Power was maximized when the bi-allelic series is balanced (4v4; 35/127 possible allelic series) and minimized when imbalanced (7v1; 8/127 possible allelic series). Uniform sampling of bi-allelic series, the approach in the more general simulations described earlier, slightly reduced power relative to balanced 4v4 allelic series due to averaging over many cases of balance and some cases of extreme imbalance. These latter, more focused simulations highlight the extent that the reduction in QTL effect size, and thus mapping power, when simulating based on Definition B, is highly dependent on the allelic series. This could be of particular importance when considering QTL that result from a causal variant inherited from a wild-derived founder, such as CAST, which will present as both imbalanced and bi-allelic.

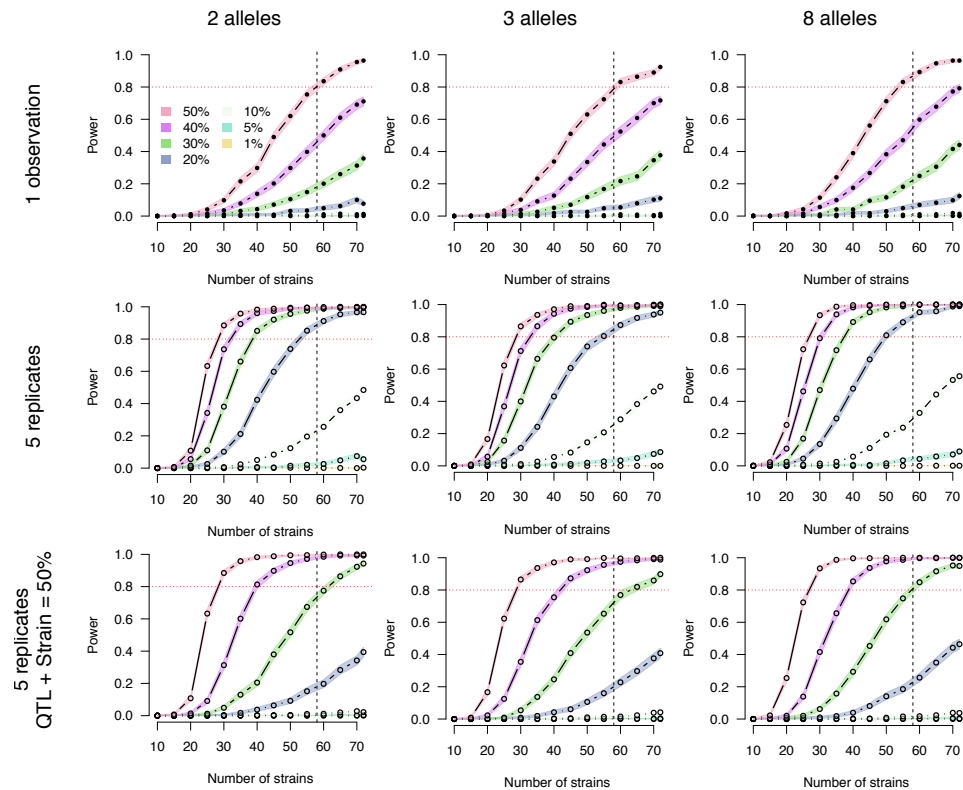
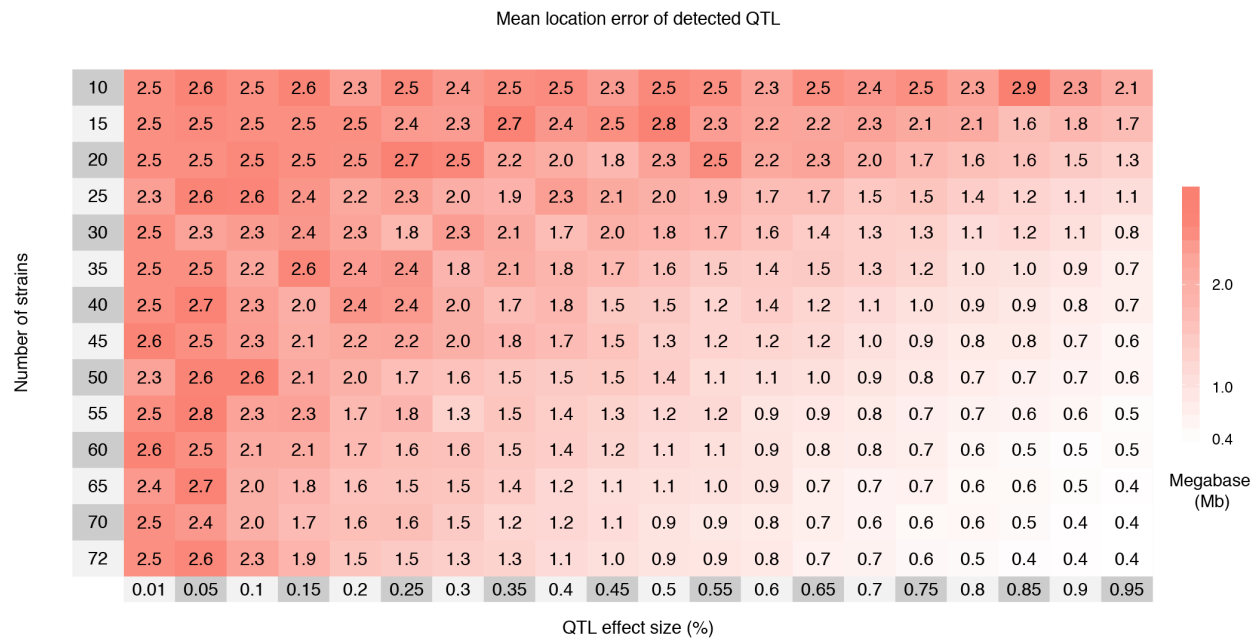


Figure 3 Power curves by number of CC strains. Results are stratified by a number of replicates, background strain effect size, and the number of functional alleles. The **[top]** row is based on a single observation per strain and no background strain effect. The **[middle]** row corresponds to five replicates per strain and no background strain effect. For the **[bottom row]**, five replicates are observed and the QTL effect size and background strain effect size sum to 50%, thus penalizing smaller QTL more harshly. The horizontal red dotted line marks 80% power. The vertical black dashed line marks 58 strains, which is currently the number of unrelated strains available from UNC. The columns, left to right, correspond to two, three, and eight functional alleles. Closed circles represent power estimates that were directly assessed, whereas open circles were interpolated. Simulations are based on Definition B.

A



B

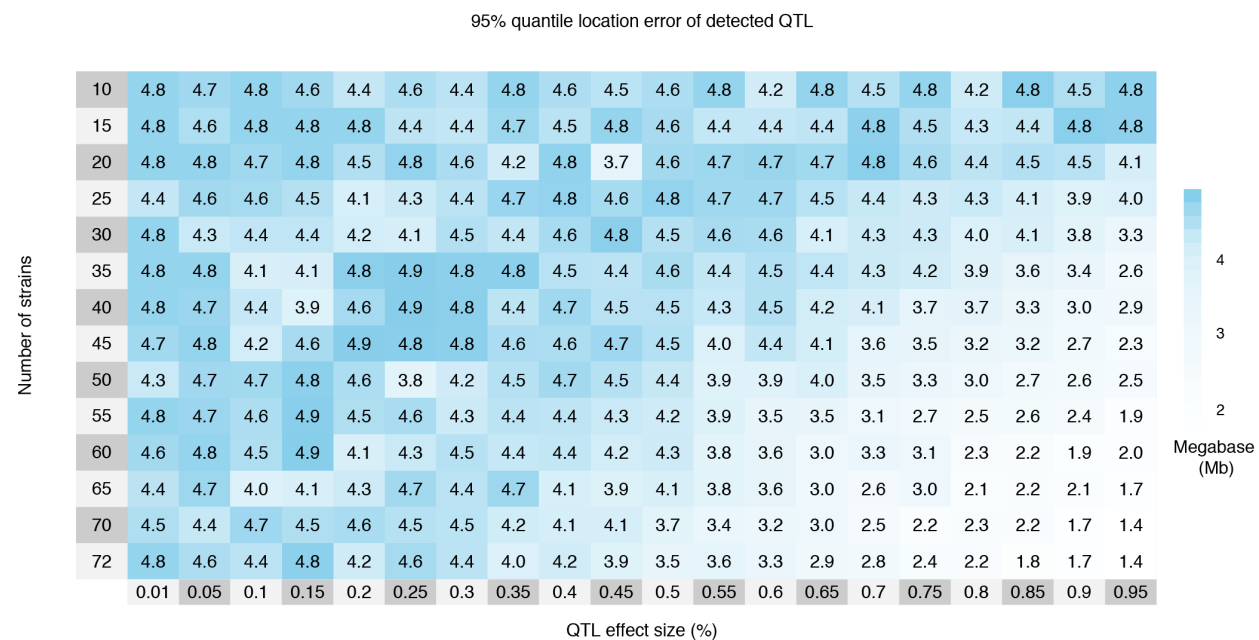


Figure 4 The mean (A) and 95% quantile (B) of location error, the distance in Mb between the detected and simulated QTL, by effect size and number of strains for 1,000 simulations of each setting. The simulations are based on Definition DAMB with an eight allele QTL, and only a single observation per strain. Cells are colored red to white with decreasing mean and blue to white with decreasing 95% quantile. Regularization of the means and 95% quantile was accomplished through averaging the observed results with pseudo-counts; see **Figure S3** for the raw measurements. Increasing the number of strains reduces the location error, both in terms of the mean and 95% quantile, more so than QTL effect size, also shown in **Figure S5**. The maximum possible location error was 5 Mb due to the 10 Mb window centered around the true QTL position used for detecting QTL.

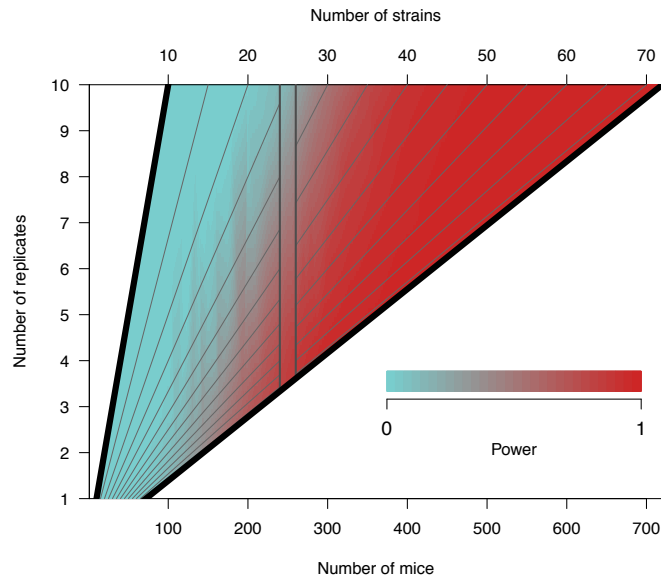


Figure 5 Heatmap of QTL mapping power by number of replicates and total number of mice in the experiment. Power is based on a QTL effect size of 20%, no background strain effect, and two functional alleles, though varying these parameters does not affect the dynamic between number of strains and replicates. The gray diagonal lines represent fixed values of the number of CC strains, ranging from 10 to 70 in intervals of five. Holding the total number of mice fixed, power is reduced as the percentage of the sample that are replicates is increased. This is illustrated with a cutout band centered on 250 mice, where power is lower at the top of the band when replicate mice are a relatively higher proportion of the total number of mice.

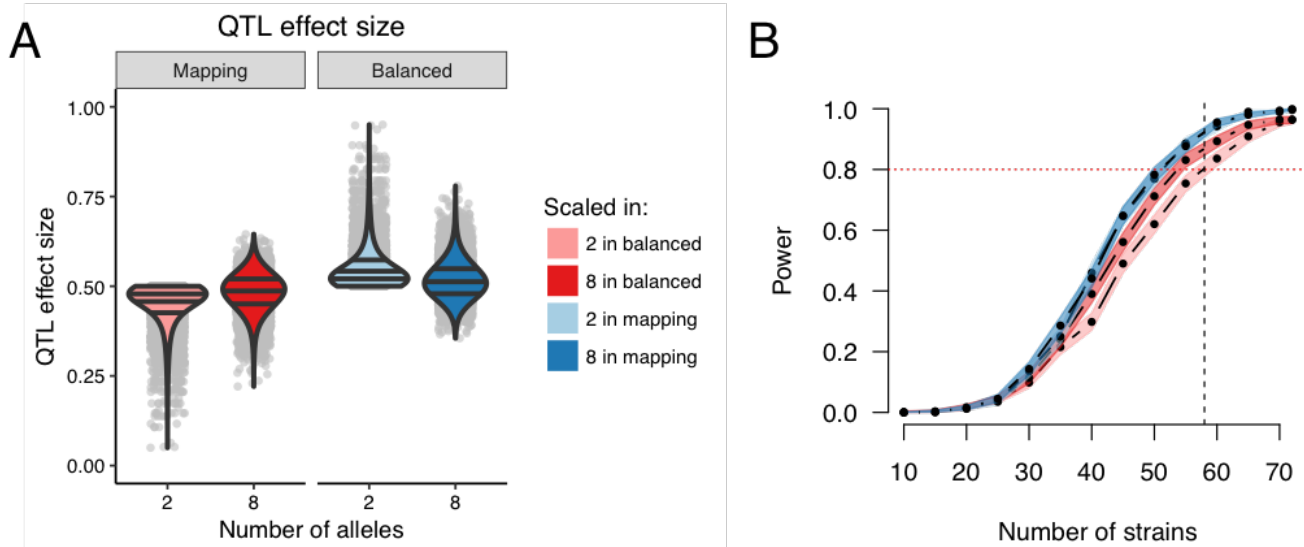


Figure 6 QTL effect sizes are in reference to a population, though effect size in the specific mapping population will determine the mapping power. Consider two populations as examples: the mapping population (definition DAMB) and a population balanced in the functional alleles (definition B). (A) QTL effect size distributions based on 10,000 simulations of the QTL for 72 strains. Using definition B, the effect sizes for the mapping population for two alleles is pink and eight alleles is red. Using definition DAMB, the effect sizes in the balanced population for two alleles is light blue and eight alleles is dark blue. Horizontal lines within the violin plots represent the 25th, 50th, and 75th quantiles from the estimated densities. Gray dots represent actual data points. (B) Power curves corresponding to the previously described settings of alleles and QTL effect size definitions. Power curves are estimated from 1,000 simulations per number of strains for a 50% QTL, no background strain effect, and a single observation per strain. The horizontal red dotted line marks 80% power. The vertical black dashed line marks 58 strains, which is currently the number of unrelated strains available from UNC.

Table 1 QTL mapping power in the Collaborative Cross based on QTL effect sizes in a balanced population (Definition B)

QTL			Power									
			30 strains			50 strains			72 strains			
1 obs ^a	3 rep ^b	5 rep ^b	2 alleles	3 alleles	8 alleles	2 alleles	3 alleles	8 alleles	2 alleles	3 alleles	8 alleles	
0.01	0.003	0.002	0.001	0.000	0.000	0.000	0.001	0.001	0.001	0.001	0.000	0.000
0.05	0.017	0.010	0.001	0.001	0.002	0.004	0.000	0.001	0.001	0.007	0.000	0.003
0.1	0.036	0.022	0.001	0.001	0.001	0.006	0.003	0.004	0.004	0.013	0.013	0.014
0.15	0.056	0.034	0.001	0.003	0.002	0.009	0.011	0.014	0.014	0.035	0.054	0.041
0.2	0.077	0.048	0.006	0.009	0.003	0.032	0.026	0.030	0.030	0.077	0.110	0.124
0.25	0.100	0.062	0.002	0.011	0.015	0.076	0.061	0.066	0.066	0.207	0.231	0.252
0.3	0.125	0.079	0.011	0.014	0.010	0.105	0.118	0.116	0.116	0.357	0.377	0.441
0.35	0.152	0.097	0.018	0.024	0.034	0.194	0.207	0.261	0.261	0.553	0.564	0.633
0.4	0.182	0.118	0.035	0.038	0.056	0.298	0.335	0.383	0.383	0.711	0.717	0.792
0.45	0.214	0.141	0.048	0.063	0.078	0.456	0.467	0.539	0.539	0.858	0.857	0.905
0.5	0.250	0.167	0.098	0.102	0.114	0.620	0.630	0.712	0.712	0.964	0.924	0.964
0.55	0.289	0.196	0.156	0.180	0.208	0.789	0.784	0.860	0.860	0.977	0.961	0.993
0.6	0.333	0.231	0.272	0.251	0.304	0.914	0.896	0.935	0.935	0.990	0.984	0.998
0.65	0.382	0.271	0.387	0.412	0.486	0.953	0.934	0.985	0.985	0.993	0.992	0.999
0.7	0.438	0.318	0.603	0.582	0.635	0.983	0.965	0.994	0.994	0.998	0.993	1.000
0.75	0.500	0.375	0.780	0.746	0.818	0.990	0.986	0.999	0.999	0.998	0.999	1.000
0.8	0.571	0.444	0.890	0.851	0.923	0.995	0.991	1.000	1.000	0.999	1.000	1.000
0.85	0.654	0.531	0.932	0.927	0.983	0.997	0.995	0.999	0.999	1.000	1.000	1.000
0.9	0.750	0.643	0.970	0.955	0.994	0.999	0.999	1.000	1.000	1.000	0.999	1.000
0.95	0.864	0.792	0.976	0.966	1.000	0.999	0.998	1.000	1.000	1.000	1.000	1.000

^a Convert QTL effect sizes from experiments with replicates to mean scale with Eq 4.

^b Based on no background strain effect.

Discussion

Now that the CC strains have been largely finalized, it is possible to investigate more deeply how, in potential mapping experiments, power is affected by factors such as the number of strains, the number of replicates, and the allelic series at the QTL. We find that the CC can powerfully map large effect QTL ($\geq 50\%$) with single observations of > 60 strains. Through the use of replicates, the power to map QTL can be greatly improved, potentially mapping QTL $\geq 20\%$ in 60 strains with 5 replicates per strain with no background strain effect. To guide the design of new CC experiments, we provide broad power curves and tables in Figure 3 and Tables 1 and S1.

The power calculations described here take advantage of realized CC genomes, allowing the power estimates to be highly specific to the available strains but also necessarily restricting the number that can be used. This differs from the simulations of Valdar *et al.* (2006a), which primarily focused on comparing potential breeding designs with numbers of strains that far exceed (500-1,000) the realized population (50-70). As such, directly comparing these studies is challenging. The closest comparison case is for a 5% QTL with 45% background strain effect with 100 simulated strains with 10 replicates, for which Valdar *et al.*

(2006a) estimates 4% power. Matching those settings with the exception of 72 strains instead of 100, and using the DAMB definition of QTL effect size, we find 0.4% power. The relatively lower power with the realized data likely reflects both reduction in the number of strains by 28% (72 to 100) and the deviations from an ideally-randomized population, such as the observed reduction in contributions from the CAST and PWK founders (Srivastava *et al.* 2017). This emphasizes the challenge in projecting the results from Valdar *et al.* (2006a) into the realized population for the purpose of designing an experiment.

We did not attempt power simulations with epistatic QTL or phenotypes with large background strain effect. From the results of Valdar *et al.* (2006a), it was clear that mapping studies in the realized CC, even with replicates, would not be well-powered in those contexts. Nonetheless, despite the reduced number of strains of realized population, we found that successful mapping experiments can be designed in the realized CC, particularly by harnessing the ability of genetic replicates to reduce random noise, as well as within the context of molecular phenotypes such as gene expression for which the genetic architecture is relatively simple.

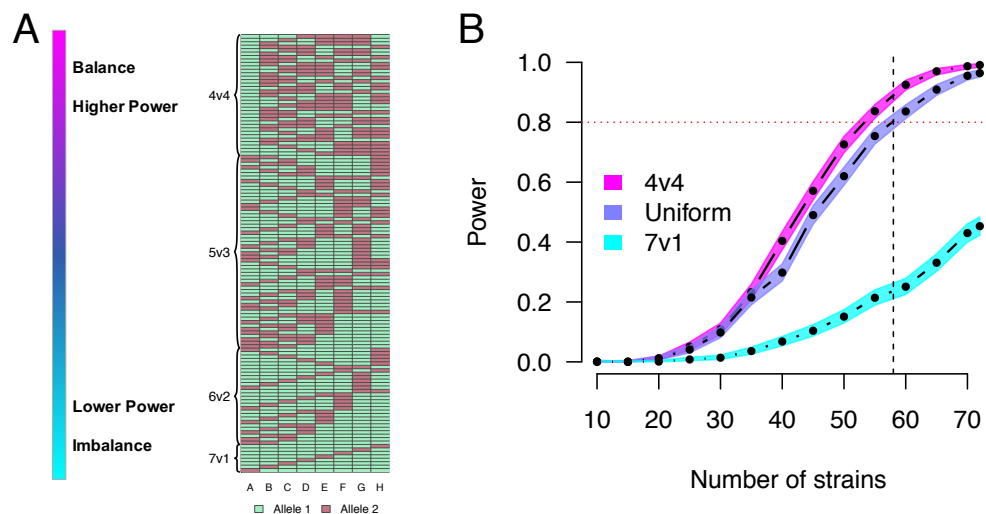


Figure 7 The balance of the allelic series for QTL with two functional alleles, and its effect on QTL mapping power. (A) The 127 possible allelic series for a bi-allelic QTL, categorized by the balance in the distribution of alleles among the CC founder strains, and ordered with balanced allelic series at the top and imbalanced at the bottom. (B) Power curves comparing three different sampling approaches for the allelic series with two functional alleles, for populations simulated to have a QTL effect size of 50% in a balanced theoretical population, with a single observation per CC strain. The horizontal red dotted line marks 80% power. The vertical black dashed line marks 58 strains, which is currently the number of unrelated strains available from UNC.

More strains vs more replicates

When holding the total number of mice fixed, we found that adding more strains improves power and reduces location error to a greater degree than does adding more replicates. Moreover, this inference was made in the absence of a background strain effect; given that replicates reduce individual-level variance but not strain-level variance, the presence of background effects would reduce the relative value of replicates yet further. These observations are consistent with the results of Valdar *et al.* (2006a) and established theoretical arguments (Soller and Beckmann 1990; Knapp and Bridges 1990).

Nonetheless, for many CC mapping experiments we predict that adding replicates will provide considerable value. First, for all but the most highly polygenic traits, mapping on the means of replicates, a strategy originally termed “replicated progeny” (Cowen 1988) or “progeny testing” (Lander and Botstein 1989), will always provide additional power. Indeed, with a limited number of strains available, and the possibility that all available strains are used, replication may sometimes be the only way power can be further increased (Belknap 1998).

Second, replicates provide not only an insurance policy against phenotyping errors, but also a way to average over batches and similar nuisance parameters (Cowen 1988), thus protecting against the negative consequences of gene by environment interactions while also providing the opportunity for such interactions to be detected [*e.g.*, Kafkafi *et al.* (2005, 2018)].

Third, replicates enable deeper phenotypic characterization and in particular measurement of strain-level phenotypes that are necessarily a function of multiple individuals. For example, treatment response phenotypes (*e.g.*, response to drug) are ideally defined in terms of counterfactual-like observations of drug-treated and vehicle-treated strain replicates [*e.g.*, Festing (2010); Crowley *et al.* (2014)] and recombinant inbred lines such

as the CC are uniquely able to combine such definitions with QTL mapping [Mosedale *et al.* (2017), but also see examples in the DSPR: Kislukhin *et al.* (2013); Najarro *et al.* (2015)]. Similarly, strain-specific phenotypic variance ideally requires replicates (Rönnegård and Valdar 2011; Ayroles *et al.* 2015). We did not consider such elaborations here, but we expect the trade-off between number of strains vs replicates may be more subtle in such cases.

Allelic series, and use of an eight allele mapping model

We found that the allelic series can strongly affect power through its influence on observed allele frequencies. Specifically, imbalanced bi-allelic QTL have significantly reduced mapping power whereas highly multi-allelic QTL do not because the potential for imbalance is reduced.

Regardless of the true allelic series at a QTL, which is unknown in practice, our statistical procedure assumed an eight allele model. For QTL with fewer functional alleles than founder strains, this assumption could reduce power due to the estimation of redundant allele effect parameters. QTL with simpler allelic series could be mapped more powerfully via a simpler QTL model, as has been seen in some MPP studies using (biallelic) SNP association (Baud *et al.* 2013; Keele *et al.* 2018). Nonetheless, incorporating simplified but unknown allelic series and their corresponding uncertainty proves challenging in practice, and the development of alternative mapping strategies that specifically account for the allelic series remains to be adequately addressed. Such approaches would likely be computationally expensive and poorly suited to simulation-based power analyses. As a standard approach, however, the eight allele model, although potentially redundant, encompasses all possible simpler allelic series, and implicitly models local epistasis and LD in a way not possible with SNP or variant association.

Use of extinct CC strains in simulations

Our simulations included genomes from CC strains that are now extinct, and also did not include all the CC strains that are currently available. This discrepancy reflects the inherent challenge of maintaining a stable genetic population resource. RI panels, such as the CC, are an approximation to an ideal: they attempt to provide reproducible genomes that can be observed multiple times as well as across multiple studies; yet, as a biological population, the genomes are mutable, and through time will accumulate mutations and drift, and even potentially go extinct.

Although the inclusion of genomes of extinct strains, or those that have drifted since the strains were genotyped, result in power calculations that do not perfectly correspond to the current CC population, they are preferable to simulated genomes, since they represent genomes that were viable at some point. We view the use of extinct genomes as realistic observations of possible genomes that reflect both the potential that more strains will become extinct or be gained from other breeding sites with time, and thus can be reasonably extended to the realized population, now and into the future.

Future use and directions

Any analysis of power is subject to the assumptions underlying that analysis. One of the advantages of simulation is the ability to evaluate the impact of many of these assumptions, as well as the consideration of new scenarios by re-running the simulation under different settings, or by elaborating the simulation itself. We have attempted to make re-running the simulations under different settings straightforward for other researchers by developing a software package for this purpose. This package could be used to investigate highly-specialized questions, such as the power for specific combinations of CC strains or assessing how the power to detect QTL varies depending on genomic position. In future work, the simulation code itself could be expanded to investigate additional topics of interest, such as how variance heterogeneity or model mis-specification influence power.

Conclusion

We used a focused simulation approach that incorporates realized CC genomes to provide more accurate estimates of QTL mapping power than were previously possible. As such, the results of our simulations provide tailored power calculations to aide the design of future QTL mapping experiments using the CC. Additionally, we evaluate how the balance of alleles at the QTL can strongly influence power to map QTL in the CC. We make available the R package SPARCC that we developed for running these simulations and analyses. It leverages an efficient model fitting approach in order to explore power in a level of detail that has previously been impractical, it is replicable, and it can be extended to user-specified questions of interest.

Acknowledgments

This work was primarily supported by the National Institute of General Medical Sciences under awards R01-GM104125 and R35-GM127000 (to W.V) and the National Institute of Environmental Health Sciences under award R01-ES024965 (to S.N.P.K). Computing resources were generously provided by the University of North Carolina Information Technology Services.

Literature Cited

- Aylor, D. L., W. Valdar, W. Foulds-mathes, R. J. Buus, R. A. Verdugo, *et al.*, 2011 Genetic analysis of complex traits in the emerging Collaborative Cross. *Genome research* **21**: 1213–22.
- Ayroles, J. F., S. M. Buchanan, C. O’Leary, K. Skutt-Kakaria, J. K. Grenier, *et al.*, 2015 Behavioral idiosyncrasy reveals genetic control of phenotypic variability. *Proceedings of the National Academy of Sciences* **112**: 6706–6711.
- Baud, A., R. Hermsen, V. Guryev, P. Stridh, D. Graham, *et al.*, 2013 Combined sequence-based and genetic mapping analysis of complex traits in outbred rats. *Nature genetics* **45**: 767–75.
- Belknap, J. K., 1998 Effect of within-strain sample size on GTL detection and mapping using recombinant inbred mouse strains. *Behavior Genetics* **28**: 29–38.
- Belknap, J. K., S. R. Mitchell, L. A. O’Toole, M. L. Helms, and J. C. Crabbe, 1996 Type I and type II error rates for quantitative trait loci (QTL) mapping studies using recombinant inbred mouse strains. *Behavior genetics* **26**: 149–60.
- Bouchet, S., M. O. Olatoye, S. R. Marla, R. Perumal, T. Tesso, *et al.*, 2017 Increased Power To Dissect Adaptive Traits in Global Sorghum Diversity Using a Nested Association Mapping Population. *Genetics* **206**: 573–585.
- Brown, L. D., T. T. Cai, and A. DasGupta, 2001 Interval Estimation for a Binomial Proportion. *Statistical Science* **16**: 101–117.
- Chesler, E. J., D. R. Miller, L. R. Branstetter, L. D. Galloway, B. L. Jackson, *et al.*, 2008 The Collaborative Cross at Oak Ridge National Laboratory: developing a powerful resource for systems genetics. *Mammalian Genome* **19**: 382–389.
- Churchill, G. A., D. C. Airey, H. Allayee, J. M. Angel, A. D. Attie, *et al.*, 2004 The Collaborative Cross, a community resource for the genetic analysis of complex traits. *Nature Genetics* **36**: 1133–1137.
- Collaborative Cross Consortium, 2012 The genome architecture of the Collaborative Cross mouse genetic reference population. *Genetics* **190**: 389–401.
- Cowen, N. M., 1988 The use of replicated progenies in marker-based mapping of QTL’s. *Theoretical and Applied Genetics* **75**: 857–862.
- Crowley, J. J., Y. Kim, A. B. Lenarcic, C. R. Quackenbush, C. J. Barrick, *et al.*, 2014 Genetics of adverse reactions to haloperidol in a mouse diallel: a drug-placebo experiment and Bayesian causal analysis. *Genetics* **196**: 321–47.
- Dell’Acqua, M., D. M. Gatti, G. Pea, F. Cattonaro, F. Coppens, *et al.*, 2015 Genetic properties of the MAGIC maize population: a new platform for high definition QTL mapping in *Zea mays*. *Genome biology* **16**: 167.
- Doerge, R. and G. Churchill, 1996 Permutation tests for multiple loci affecting a quantitative character. *Genetics* **142**: 285–94.
- Donoghue, L. J., A. Livraghi-Buttrico, K. M. McFadden, J. M. Thomas, G. Chen, *et al.*, 2017 Identification of trans Protein QTL for Secreted Airway Mucins in Mice and a Causal Role for Bpifb1. *Genetics* **207**: 801–812.
- Dudbridge, F. and B. P. Koeleman, 2004 Efficient Computation of Significance Levels for Multiple Associations in Large Studies of Correlated Data, Including Genomewide Association Studies. *The American Journal of Human Genetics* **75**: 424–435.
- Durrant, C., H. Tayem, B. Yalcin, J. Cleak, L. Goodstadt, *et al.*, 2011 Collaborative Cross mice and their power to map host susceptibility to *Aspergillus fumigatus* infection. *Genome research* **21**: 1239–48.
- Falke, K. C. and M. Frisch, 2011 Power and false-positive rate in QTL detection with near-isogenic line libraries. *Heredity* **106**:

- 576–584.
- Ferris, M. T., D. L. Aylor, D. Bottomly, A. C. Whitmore, L. D. Aicher, *et al.*, 2013 Modeling Host Genetic Regulation of Influenza Pathogenesis in the Collaborative Cross. *PLoS Pathogens* **9**: e1003196.
- Festing, M. F. W., 2010 Inbred strains should replace outbred stocks in toxicology, safety testing, and drug development. *Toxicologic Pathology* **38**: 681–690.
- Fu, C.-P., C. E. Welsh, F. P.-M. de Villena, and L. McMillan, 2012 Inferring ancestry in admixed populations using microarray probe intensities. In *Proceedings of the ACM Conference on Bioinformatics, Computational Biology and Biomedicine - BCB '12*, pp. 105–112, New York, New York, USA, ACM Press.
- Gatti, D. M., K. L. Svenson, A. Shabalin, L.-Y. Wu, W. Valdar, *et al.*, 2014 Quantitative Trait Locus Mapping Methods for Diversity Outbred Mice. *G3 (Bethesda, Md.)* **4**: 1623–1633.
- Graham, J. B., J. L. Swarts, M. Mooney, G. Choonoo, S. Jeng, *et al.*, 2017 Extensive Homeostatic T Cell Phenotypic Variation within the Collaborative Cross. *Cell reports* **21**: 2313–2325.
- Gralinski, L. E., M. T. Ferris, D. L. Aylor, A. C. Whitmore, R. Green, *et al.*, 2015 Genome Wide Identification of SARS-CoV Susceptibility Loci Using the Collaborative Cross. *PLoS genetics* **11**: e1005504.
- Haley, C. S. and S. A. Knott, 1992 A simple regression method for mapping quantitative trait loci in line crosses using flanking markers. *Heredity* **69**: 315–24.
- Kaeppler, S. M., 1997 Quantitative trait locus mapping using sets of near-isogenic lines: Relative power comparisons and technical considerations. *Theoretical and Applied Genetics* **95**: 384–392.
- Kafkafi, N., J. Agassi, E. J. Chesler, J. C. Crabbe, W. E. Crusio, *et al.*, 2018 Reproducibility and replicability of rodent phenotyping in preclinical studies. *Neuroscience and Biobehavioral Reviews* **87**: 218–232.
- Kafkafi, N., Y. Benjamini, A. Sakov, G. I. Elmer, and I. Golani, 2005 Genotype-environment interactions in mouse behavior: a way out of the problem. *Proceedings of the National Academy of Sciences of the United States of America* **102**: 4619–24.
- Keele, G. R., J. W. Prokop, H. He, K. Holl, J. Littrell, *et al.*, 2018 Genetic Fine-Mapping and Identification of Candidate Genes and Variants for Adiposity Traits in Outbred Rats. *Obesity* **26**: 213–222.
- Kelada, S. N. P., 2016 Plethysmography Phenotype QTL in Mice Before and After Allergen Sensitization and Challenge. *G3 (Bethesda, Md.)* **6**: 2857–2865.
- Kelada, S. N. P., D. L. Aylor, B. C. E. Peck, J. F. Ryan, U. Tavarez, *et al.*, 2012 Genetic Analysis of Hematological Parameters in Incipient Lines of the Collaborative Cross. *G3 (Bethesda, Md.)* **2**: 157–165.
- King, E. G. and A. D. Long, 2017 The Beavis Effect in Next-Generation Mapping Panels in *Drosophila melanogaster*. *G3* **7**: 1643 LP – 1652.
- King, E. G., S. J. Macdonald, and A. D. Long, 2012 Properties and power of the *Drosophila* synthetic population resource for the routine dissection of complex traits. *Genetics* **191**: 935–949.
- Kislukhin, G., E. G. King, K. N. Walters, S. J. Macdonald, and A. D. Long, 2013 The Genetic Architecture of Methotrexate Toxicity Is Similar in *Drosophila melanogaster* and Humans. *G3: Genes, Genomes, Genetics* **3**: 1301–1310.
- Klasen, J. R., H. P. Piepho, and B. Stich, 2012 QTL detection power of multi-parental RIL populations in *Arabidopsis thaliana*. *Heredity* **108**: 626–632.
- Knapp, S. J. and W. C. Bridges, 1990 Using molecular markers to estimate quantitative trait locus parameters: power and genetic variances for unreplicated and replicated progeny. *Genetics* **126**: 769–77.
- Kover, P. X., W. Valdar, J. Trakalo, N. Scarcelli, I. M. Ehrenreich, *et al.*, 2009 A Multiparent Advanced Generation Inter-Cross to fine-map quantitative traits in *Arabidopsis thaliana*. *PLoS genetics* **5**: e1000551.
- Lander, E. S. and D. Botstein, 1989 Mapping mendelian factors underlying quantitative traits using RFLP linkage maps. *Genetics* **121**: 185–99.
- Li, H., P. Bradbury, E. Ersoz, E. S. Buckler, and J. Wang, 2011 Joint QTL linkage mapping for multiple-cross mating design sharing one common parent. *PloS one* **6**: e17573.
- Liu, E. Y., Q. Zhang, L. McMillan, F. P.-M. de Villena, and W. Wang, 2010 Efficient genome ancestry inference in complex pedigrees with inbreeding. *Bioinformatics* **26**: i199–i207.
- Lorè, N. I., F. A. Iraqi, and A. Bragonzi, 2015 Host genetic diversity influences the severity of *Pseudomonas aeruginosa* pneumonia in the Collaborative Cross mice. *BMC genetics* **16**: 106.
- Mackay, T. F. C., S. Richards, E. A. Stone, A. Barbadilla, J. F. Ayroles, *et al.*, 2012 The *Drosophila melanogaster* Genetic Reference Panel. *Nature* **482**: 173–8.
- Martínez, O. and R. N. Curnow, 1992 Estimating the locations and the sizes of the effects of quantitative trait loci using flanking markers. *Theor. Appl. Genet.* **85**: 480–488.
- Mathes, W. F., D. L. Aylor, D. R. Miller, G. A. Churchill, E. J. Chesler, *et al.*, 2011 Architecture of energy balance traits in emerging lines of the Collaborative Cross. *American Journal of Physiology-Endocrinology and Metabolism* **300**: E1124–E1134.
- Morgan, A. P., C.-P. Fu, C.-Y. Kao, C. E. Welsh, J. P. Didion, *et al.*, 2016 The Mouse Universal Genotyping Array: From Substrains to Subspecies. *G3: Genes, Genomes, Genetics* **6**: 263–279.
- Mosedale, M., Y. Kim, W. J. Brock, S. E. Roth, T. Wiltshire, *et al.*, 2017 Candidate Risk Factors and Mechanisms for Tolvaptan-Induced Liver Injury Are Identified Using a Collaborative Cross Approach. *Toxicological Sciences* **156**: kfw269.
- Mott, R., C. J. Talbot, M. G. Turri, A. C. Collins, and J. Flint, 2000 A method for fine mapping quantitative trait loci in outbred animal stocks. *PNAS* **97**: 12649–54.
- Najarro, M. A., J. L. Hackett, B. R. Smith, C. A. Highfill, E. G. King, *et al.*, 2015 Identifying Loci Contributing to Natural Variation in Xenobiotic Resistance in *Drosophila*. *PLoS Genetics* **11**: 1–25.
- Noble, L. M., I. Chelo, T. Guzella, B. Afonso, D. D. Riccardi, *et al.*, 2017 Polygenicity and Epistasis Underlie Fitness-Proximal Traits in the *Caenorhabditis elegans* Multiparental Experimental Evolution (CeMEE) Panel. *Genetics* **207**: genetics.300406.2017.
- Orgel, K., J. M. Meekens, P. Ye, L. Fotsch, R. Guo, *et al.*, 2018 Genetic Diversity Between Mouse Strains Allows Identification of CC027/GeniUnc as an Orally Reactive Model of Peanut Allergy. *The Journal of allergy and clinical immunology*.
- Peirce, J. L., L. Lu, J. Gu, L. M. Silver, and R. W. Williams, 2004 A new set of BXD recombinant inbred lines from advanced intercross populations in mice. *BMC genetics* **5**: 7.
- Pfaff, B. and A. McNeil, 2018 *evir: Extreme Values in R*. R package version 1.7-4.
- Philip, V. M., G. Sokoloff, C. L. Ackert-Bicknell, M. Striz,

- L. Branstetter, *et al.*, 2011 Genetic analysis in the Collaborative Cross breeding population. *Genome Research* **21**: 1223–1238.
- Phillippi, J., Y. Xie, D. R. Miller, T. A. Bell, Z. Zhang, *et al.*, 2014 Using the emerging Collaborative Cross to probe the immune system. *Genes & Immunity* **15**: 38–46.
- R Core Team, 2018 *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Ram, R., M. Mehta, L. Balmer, D. M. Gatti, and G. Morahan, 2014 Rapid identification of major-effect genes using the collaborative cross. *Genetics* **198**: 75–86.
- Rasmussen, A. L., A. Okumura, M. T. Ferris, R. Green, F. Feldmann, *et al.*, 2014 Host genetic diversity enables Ebola hemorrhagic fever pathogenesis and resistance. *Science (New York, N.Y.)* **346**: 987–91.
- Rogala, A. R., A. P. Morgan, A. M. Christensen, T. J. Gooch, T. A. Bell, *et al.*, 2014 The Collaborative Cross as a resource for modeling human disease: CC011/Unc, a new mouse model for spontaneous colitis. *Mammalian genome* **25**: 95–108.
- Rönnegård, L. and W. Valdar, 2011 Detecting major genetic loci controlling phenotypic variability in experimental crosses. *Genetics* **188**: 435–447.
- Rutledge, H., D. L. Aylor, D. E. Carpenter, B. C. Peck, P. Chines, *et al.*, 2014 Genetic regulation of Zfp30, CXCL1, and neutrophilic inflammation in murine lung. *Genetics* **198**: 735–745.
- Shorter, J. R., F. Odet, D. L. Aylor, W. Pan, C.-Y. Kao, *et al.*, 2017 Male Infertility Is Responsible for Nearly Half of the Extinction Observed in the Mouse Collaborative Cross. *Genetics* **206**: 557–572.
- Shusterman, A., Y. Salyma, A. Nashef, M. Soller, A. Wilensky, *et al.*, 2013 Genotype is an important determinant factor of host susceptibility to periodontitis in the Collaborative Cross and inbred mouse populations. *BMC genetics* **14**: 68.
- Soller, M. and J. S. Beckmann, 1990 Marker-based mapping of quantitative trait loci using replicated progenies. *Theoretical and Applied Genetics* **80**: 205–208.
- Srivastava, A., A. P. Morgan, M. L. Najarian, V. K. Sarsani, J. S. Sigmon, *et al.*, 2017 Genomes of the mouse Collaborative Cross. *Genetics* **206**: 537–556.
- Svenson, K. L., D. M. Gatti, W. Valdar, C. E. Welsh, R. Cheng, *et al.*, 2012 High-resolution genetic mapping using the Mouse Diversity outbred population. *Genetics* **190**: 437–47.
- Takuno, S., R. Terauchi, and H. Innan, 2012 The power of QTL mapping with RILs. *PloS one* **7**: e46545.
- Threadgill, D. W. and G. A. Churchill, 2012 Ten Years of the Collaborative Cross. *Genetics* **190**: 291–294.
- Threadgill, D. W., K. W. Hunter, and R. W. Williams, 2002 Genetic dissection of complex and quantitative traits: from fantasy to reality via a community effort. *Mammalian genome : official journal of the International Mammalian Genome Society* **13**: 175–8.
- Valdar, W., J. Flint, and R. Mott, 2006a Simulating the Collaborative Cross: power of quantitative trait loci detection and mapping resolution in large sets of recombinant inbred strains of mice. *Genetics* **172**: 1783–97.
- Valdar, W., L. C. Solberg, D. Gauguier, S. Burnett, P. Klenerman, *et al.*, 2006b Genome-wide genetic association of complex traits in heterogeneous stock mice. *Nature Genetics* **38**: 879–887.
- Venables, W. N. and B. D. Ripley, 2002 *Modern Applied Statistics with S*. Springer, New York, fourth edition, ISBN 0-387-95457-0.
- Venkatratnam, A., S. Furuya, O. Kosyik, A. Gold, W. Bodnar, *et al.*, 2017 Collaborative Cross Mouse Population Enables Refinements to Characterization of the Variability in Toxicokinetics of Trichloroethylene and Provides Genetic Evidence for the Role of PPAR Pathway in Its Oxidative Metabolism. *Toxicological Sciences* **158**: 48–62.
- Vered, K., C. Durrant, R. Mott, and F. A. Iraqi, 2014 Susceptibility to Klebsiella pneumoniae infection in collaborative cross mice is a complex trait controlled by at least three loci acting at different time points. *BMC genomics* **15**: 865.
- Welsh, C. E., D. R. Miller, K. F. Manly, J. Wang, L. McMillan, *et al.*, 2012 Status and access to the Collaborative Cross population. *Mammalian genome : official journal of the International Mammalian Genome Society* **23**: 706–12.
- Yalcin, B., J. Flint, and R. Mott, 2005 Using progenitor strain information to identify quantitative trait nucleotides in outbred mice. *Genetics* **171**: 673–81.
- Yamamoto, E., H. Iwata, T. Tanabata, R. Mizobuchi, J.-i. Yonemaru, *et al.*, 2014 Effect of advanced intercrossing on genome structure and on the power to detect linked quantitative trait loci in a multi-parent population: a simulation study in rice. *BMC genetics* **15**: 50.
- Yang, H., J. R. Wang, J. P. Didion, R. J. Buus, T. A. Bell, *et al.*, 2011 Subspecific origin and haplotype diversity in the laboratory mouse. *Nature genetics* **43**: 648–55.
- Yu, J., J. B. Holland, M. D. McMullen, and E. S. Buckler, 2008 Genetic design and statistical power of nested association mapping in maize. *Genetics* **178**: 539–551.
- Zhang, Z., W. Wang, and W. Valdar, 2014 Bayesian modeling of haplotype effects in multiparent populations. *Genetics* **198**: 139–56.
- Zheng, C., M. P. Boer, and F. A. van Eeuwijk, 2015 Reconstruction of Genome Ancestry Blocks in Multiparental Populations. *Genetics* **200**: 1073–1087.
- Zou, F., Z. Xu, and T. Vision, 2006 Assessing the significance of quantitative trait loci in replicable mapping populations. *Genetics* **174**: 1063–8.

Appendix A: Available options in the R package SPARCC

Phenotype simulation

The output of the CC data simulation is a matrix of outcomes, where each column of $\mathbf{y}^{(s)}$ is the phenotype generated by Eq 1. The matrix of phenotypes is simulated using the `sim.CC.data()` function. The following options control various aspects of the simulation and the simulated QTL.

QTL effect:

- `qtl.effect.size`
 - $0 \leq \text{qtl.effect.size} < 1 - \text{strain.effect.size}$
 - This argument represents h_{QTL}^2 , such that $\beta \sim N(\mathbf{0}, \mathbf{I}h_{\text{QTL}}^2)$.
 - A specific β can be specified with the `beta` argument, though it will be scaled to match `qtl.effect.size`. If `beta=NULL`, then β is sampled accordingly.
- `num.alleles` (DEFAULT = 8)
 - $2 \leq \text{num.alleles} \leq 8$
- `M.ID`
 - Rather than specifying `num.alleles` and then sampling \mathbf{M} , these can be fixed with the `M.ID` argument.

- Expects strings of the form "A,B,C,D,E,F,G,H", with each letter corresponding to a founder strain, taking an integer value 0-7, representing functional alleles.
- Example: M.ID="0,0,0,0,1,1,1,1" represents a bi-allelic causal variant, in which the first four strains have one allele, and the last four having the other.
- `CC.lines` (DEFAULT = NULL)
 - This argument allows the user to provide a vector of CC strain IDs on which to base the power calculation. The CC genomes, along with locus, will determine \mathbf{D} in Eq 2.
 - If `CC.lines` = NULL, then SPARCC will sample `num.lines` from all available strains.
 - * `vary.lines` (DEFAULT = TRUE)
 - If `vary.lines` = TRUE, the set of strains for each simulation will be sampled and vary.
 - If `vary.lines` = FALSE, the set of strains will be sampled once, and used for each simulated outcome.
- `locus` (DEFAULT = NULL)
 - This argument allows the user to specify a specific locus for the simulated QTL, in effect determining the haplotype dosage matrix \mathbf{D} .
 - If the argument is left empty, SPARCC will sample loci uniformly from the CC genomes, thus providing power estimates averaged over genomic positions.
- `impute` (DEFAULT = TRUE)
 - If `impute`=TRUE, then \mathbf{D} in Eq 2 is sampled from the probabilistically reconstructed diplotypes at the QTL

$$\mathbf{D}_i \sim \text{Cat}(\mathbf{P}_i) \quad (5)$$

where $\text{Cat}(\cdot)$ is a categorical distribution and \mathbf{P} is a matrix of diplotype probabilities for the sampled CC genomes at the QTL.

- If `impute`=FALSE, then $\mathbf{D} = \mathbf{P}$ in the simulation procedure.
- `scale.qtl.mode` (DEFAULT = "B")
 - If `scale.qtl.mode`="B", $V(2\beta)$ is scaled to `qtl.effect.size`, where `var` is the maximum likelihood estimate of variance rather than sample variance, such that `var` = $\frac{n-1}{n}s$, where s^2 is the sample variance and n is the number of individuals. This scaling sets the QTL effect size with respect to a theoretical population that is evenly balanced with respect to functional alleles.
 - If `scale.qtl.mode`="MB", $V(2\mathbf{M}\beta)$ is scaled to `qtl.effect.size`, setting the QTL effect size to a theoretical population with a specific frequency of functional alleles among the founder strains. The effect size here is dependent on \mathbf{M} but independent of \mathbf{D} , and the proportion of variance explained by the QTL in the mapping population will deviate from h_{QTL}^2 due to imbalance \mathbf{D} but not in \mathbf{M} .
 - If `scale.qtl.mode`="DAMB", $V(\mathbf{DAM}\beta)$ is scaled to `qtl.effect.size`, setting the QTL effect size to a specific set of CC strains and allele series.
 - If `scale.qtl.mode`="none", β is not scaled, allowing the user to specify an effect vector and not have it modified.

Strain effect:

- `strain.effect.size`
 - $0 \leq \text{strain.effect.size} \leq 1 - \text{qtl.effect.size}$
 - This argument specifies h_{strain}^2 , such that $\mathbf{u} \sim N(\mathbf{0}, \mathbf{I}h_{\text{strain}}^2)$.

Additional options:

- `num.sim`
 - This argument specifies SPARCC to simulate s samples of \mathbf{y} from Eq 1.
- `num.replicates`
 - This argument allows the user to set the number r of replicates of each CC strain. The reproducibility of CC genomes is an important feature, allowing noise variation to be reduced.
 - SPARCC currently requires all CC lines to have the same number of replicates.

Genome scans

Each phenotype from the simulation procedure is then evaluated via QTL mapping. The SPARCC function for running genome scans from the simulated data is `run.sim.scans()`, The primary argument is `sim.data`. There are additional arguments to restrict the scans to a subset of the chromosomes, to a subset of the simulated phenotypes, or to a subset of loci. Last, the user can provide the pre-computed QR decompositions and specify whether the output should return those decompositions for further use with additional simulations, although this can be expensive in terms of memory.

Appendix B: QR decomposition for fast regression

To maximize power to detect QTL while controlling the FPR, permutations to determine significance thresholds are needed, which is computationally expensive and thus the underlying regression functionality must be highly optimized. We accomplish this through the QR matrix decomposition, which we will describe briefly (Venables and Ripley 2002).

Let $\mathbf{X} = \mathbf{P}\mathbf{A}$ be the $n \times m$ design matrix included in Eq 3, with $m = 8$. The solution for β from the least squares normal equations is $\hat{\beta} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$. Through the QR decomposition, $\mathbf{X} = \mathbf{Q}\mathbf{R}$, for which \mathbf{Q} is an $n \times p$ orthonormal matrix ($\mathbf{Q}^T\mathbf{Q} = \mathbf{I}$) and \mathbf{R} is a $m \times m$ upper triangular matrix. With matrix algebra, it is fairly straightforward to show that $\hat{\beta} = \mathbf{R}^{-1}\mathbf{Q}^T\mathbf{y}$, which is also more numerically stable than calculating $\hat{\beta}$ through $(\mathbf{X}^T\mathbf{X})^{-1}$. After solving for $\hat{\beta}$, the RSS, and ultimately $\log P$, can be rapidly calculated. Because our simulation approach involves regressing many permuted outcomes $\mathbf{U}_p\mathbf{y}^{(s)}$, where \mathbf{U}_p is a permutation matrix that re-orders $\mathbf{y}^{(s)}$ randomly, on the same design matrices, computational efficiency can be vastly increased by pre-computing and saving the QR decompositions for all \mathbf{X} .

Once the QR decomposition has been stored for a design matrix \mathbf{X}_j , j indexing locus, it is highly computationally efficient to conduct additional tests for any \mathbf{y} , thus encompassing all permuted outcomes $\mathbf{U}_p\mathbf{y}$. If \mathbf{X}_j is the same across S simulations, the boost in computation can extend beyond permutations to samples of $\mathbf{y}^{(s)}$, as is the case when the set of CC strains is fixed. In effect, two cases result for our R package SPARCC: when the set of CC strains is fixed, and when the set varies.

- Fixed set of CC strains
 1. Store QR decompositions of X_j for $j = 1, 2, \dots, J$
 2. Run genome scans for $y^{(s)}$ and $U_p y^{(s)}$ for $s = 1, 2, \dots, S \times p = 1, 2, \dots, P$
- Varied set of CC strains
 1. Store QR decompositions of X_{js} for $j = 1, 2, \dots, J$
 2. Run genome scans for $y^{(s)}$ and $U_p y^{(s)}$ for $p = 1, 2, \dots, P$
 3. Repeat steps 1 and 2 for $s = 1, 2, \dots, S$

Varying the sets of CC strains increases computation time linearly with respect to S . If the investigators do not have a predefined set of strains, it is appropriate that this source of variability be incorporated into the power calculation.

Appendix C: Computing

Simple analysis code for SPARCC

This example is computationally efficient because the CC strains are not varied across simulations. Varying CC strains increases the computational expense, which was done for the reported results.

```
#####
### Useful functions for parsing haplotype data
> h <- DiploprobReader$new("./sparcc_cache/")
> set.seed(100)

### Grabbing random sample of 65 CC strains
> these.cc.lines <- sample(h$getSubjects(), size=65)

> library(sparcc)
### Simulate 5 data sets:
#### Specified 65 CC strains
#### 5 replicates of each
#### 2 functional alleles, allelic series not specified
#### QTL effect size of 10
#### Background Strain effect of 0
> simple.data <- sim.CC.data(genomecache="./sparcc_cache/",
                           CC.lines=these.cc.lines,
                           num.replicates=5,
                           num.sim=5,
                           num.alleles=2,
                           qtl.effect.size=0.1,
                           strain.effect.size=0)

### Genome scans
> simple.scans <- run.sim.scans(sim.data=simple.data,
                              return.all.sim.qr=TRUE)

### Generating permutation index
> perm.index <- generate.perm.matrix(num.lines=65,
                                    num.perm=100)

### Permutation scans
> thresh.scans <- run.perm.scans(perm.matrix=perm.index,
                               sim.CC.object=simple.data,
                               sim.CC.scans=simple.scans)

### Calculating significance thresholds
> all.thresh <- get.thresholds(thresh.scans=thresh.scans)
```

```
### Power estimate
> pull.power(sim.scans=simple.scans,
             thresh=all.thresh)
[1] 0.8

### Plot a genome scan of a single simulated phenotype
> single.sim.plot(simple.scans,
                 thresh=all.thresh,
                 phenotype.index=1)

#####
```

Run-time performance of tutorial

The simple example was run locally on an Early 2015 MacBook Pro with a 2.9 GHz Intel Core i5 processor and 8 GB of RAM. The data simulation and genome-wide scans for five phenotypes took 34.13 seconds. Computational time increases linearly with number of phenotypes simulated. Computational times will also decrease for lower numbers of CC strains. 100 permutations for 5 simulated phenotypes took 8.73 minutes. Although the time expense for these simulations is not trivial, the overall process is highly optimized; this simple example involves fitting 5 phenotypes \times 17900 loci \times 100 permutation alternative models. The process can be sped up using a parallel computing environment, as we do with the following large scale analysis. Highly specific power calculations for an experiment are feasible on a local computer using a single core.

Large-scale computing environment and performance

We performed 1,000 simulations (in batches of 100) for each combination of the parameters, resulting in 8,400 individual jobs. These jobs were submitted in parallel to a distributed computing cluster (<http://its.unc.edu/rc-services/killdevil-cluster/>). Runtime varied depending on parameter settings and the hardware used, with the longest jobs taking approximately seven hours to complete.

Power projection code for SPARCC

This example code demonstrates how to calculate power estimates from the dataset stored within SPARCC using projection and interpolation. This code produces **Figure S1**.

```
#####
### r1.dat is included in SPARCC
### Project and interpolate power estimates for
### experiments with 3 replicates
> r3.interp.dat <- interpolate.table(r1.results=r1.dat,
                                   num.replicates=3,
                                   n.alleles=2)

### Project and interpolate power estimates for
### experiments with 5 replicates
> r5.interp.dat <- interpolate.table(r1.results=r1.dat,
                                   num.replicates=5,
                                   n.alleles=2)

### Plotting power curves
> power.plot(results=r1.dat,
             qtl.effect.size=0.3,
             n.alleles=2)
> add.curve.to.power.plot(results=r3.interp.dat,
                          qtl.effect.size=0.3,
                          n.alleles=2)
> add.curve.to.power.plot(results=r5.interp.dat,
                          qtl.effect.size=0.3,
```


n.alleles=2)

See Figure S1 for corresponding plot

###-----

Appendix D: CC strains

The founder haplotype pair probability data for 72 CC strains is available at <http://csbio.unc.edu/CCstatus/index.py?run=FounderProbs>, which we used for these reported power calculations. These include the following strains:

CC001, CC002, CC003, CC004, CC005, CC006, CC007, CC008, CC009, CC010, CC011, CC012, CC013, CC014, CC015, CC016, CC017, CC018, CC019, CC020, CC021, CC022, CC023, CC024, CC025, CC026, CC027, CC028, CC029, CC030, CC031, CC032, CC033, CC034, CC035, CC036, CC037, CC038, CC039, CC040, CC041, CC042, CC043, CC044, CC045, CC046, CC047, CC048, CC049, CC050, CC051, CC052, CC053, CC054, CC055, CC056, CC057, CC058, CC059, CC060, CC061, CC062, CC063, CC065, CC068, CC070, CC071, CC072, CC073, CC074, CC075, CC076

The availability of CC strains will fluctuate with time, based on the numbers of the current colonies, as well as potential extinctions of strains or acquisition of new strains. From a personal correspondence with Darla Miller, UNC is currently planning to maintain and distribute the following 59 strains:

CC001, CC002, CC003, CC004, CC005, CC006, CC007, CC008, CC009, CC010, CC011, CC012, CC013, CC015, CC016, CC017, CC019, CC021, CC023, CC024, CC025, CC026, CC027, CC029, CC030, CC031, CC032, CC033, CC035, CC036, CC037, CC038, CC039, CC040, CC041, CC042, CC043, CC044, CC045, CC046, CC049, CC051, CC053, CC055, CC057, CC058, CC059, CC060, CC061, CC062, CC065, CC068, CC071, CC072, CC078, CC079, CC080, CC081, CC083

There are five strains (CC078, CC079, CC080, CC081, CC083) of the currently available 59 from UNC whose founder haplotype probabilities were not available on the website at the time of this work, and are thus not included in these simulations and results. The five genomes have since been added to the website. CC051 and CC059 are derived from the same breeding funnel, and thus more related than independent CC strains.

Appendix E: Additive Model and Allelic Series Matrices

Additive Matrix

	A	B	C	D	E	F	G	H
AA	2	0	0	0	0	0	0	0
BB	0	2	0	0	0	0	0	0
CC	0	0	2	0	0	0	0	0
DD	0	0	0	2	0	0	0	0
EE	0	0	0	0	2	0	0	0
FF	0	0	0	0	0	2	0	0
GG	0	0	0	0	0	0	2	0
HH	0	0	0	0	0	0	0	2
AB	1	1	0	0	0	0	0	0
AC	1	0	1	0	0	0	0	0
AD	1	0	0	1	0	0	0	0
AE	1	0	0	0	1	0	0	0
AF	1	0	0	0	0	1	0	0
AG	1	0	0	0	0	0	1	0
AH	1	0	0	0	0	0	0	1
BC	0	1	1	0	0	0	0	0
BD	0	1	0	1	0	0	0	0
BE	0	1	0	0	1	0	0	0
BF	0	1	0	0	0	1	0	0
BG	0	1	0	0	0	0	1	0
BH	0	1	0	0	0	0	0	1
CD	0	0	1	1	0	0	0	0
CE	0	0	1	0	1	0	0	0
CF	0	0	1	0	0	1	0	0
CG	0	0	1	0	0	0	1	0
CH	0	0	1	0	0	0	0	1
DE	0	0	0	1	1	0	0	0
DF	0	0	0	1	0	1	0	0
DG	0	0	0	1	0	0	1	0
DH	0	0	0	1	0	0	0	1
EF	0	0	0	0	1	1	0	0
EG	0	0	0	0	1	0	1	0
EH	0	0	0	0	1	0	0	1
FG	0	0	0	0	0	1	1	0
FH	0	0	0	0	0	1	0	1
GH	0	0	0	0	0	0	1	1

We can use matrices to specify simplifying linear combinations of the 36 diplotypes. The additive model matrix **A** is commonly used, and we use it here. Post-multiplication of the diplotype design matrix **D** with the **A** rotates the diplotypes at

the locus to dosages of the founder haplotypes. If there is no uncertainty on the diplotype identities, **DA** will be the matrix of founder haplotype counts at the locus.

Allelic series matrices

We explore the influence of the allelic series on QTL mapping power through the simulation procedure. The QTL mapping procedure estimates separate parameters for each founder, though in reality, there are likely fewer functional alleles. We denote the q^{th} functional allele as k_q . The allelic series can be sampled and encoded in the M.ID argument within the `sim.CC.data()` function of SPARCC. Below are examples of balanced (4v4) and unbalanced (7v1) bi-allelic series, as well as tri-allelic series.

Allelic series with eight alleles (maximum)

$$\begin{aligned}
 & \text{M.ID} = "0,1,2,3,4,5,6,7" \\
 & \begin{matrix} & k_0 & k_1 & k_2 & k_3 & k_4 & k_5 & k_6 & k_7 \\ \mathbf{M} = \mathbf{I} = & \begin{bmatrix} A & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ B & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ C & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ D & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ E & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ F & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ G & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ H & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \end{matrix}
 \end{aligned}$$

Example balanced (4v4) bi-allelic series

$$\begin{aligned}
 & \text{M.ID} = "0,1,0,0,1,0,1,1" \\
 & \begin{matrix} & k_0 & k_1 \\ \mathbf{M} = & \begin{bmatrix} A & 1 & 0 \\ B & 0 & 1 \\ C & 1 & 0 \\ D & 1 & 0 \\ E & 0 & 1 \\ F & 1 & 0 \\ G & 0 & 1 \\ H & 0 & 1 \end{bmatrix} \end{matrix}
 \end{aligned}$$

$$\begin{aligned}
 & \text{M.ID} = "0,1,1,1,0,0,1,0" \\
 & \begin{matrix} & k_0 & k_1 \\ \mathbf{M} = & \begin{bmatrix} A & 1 & 0 \\ B & 0 & 1 \\ C & 0 & 1 \\ D & 0 & 1 \\ E & 1 & 0 \\ F & 1 & 0 \\ G & 0 & 1 \\ H & 1 & 0 \end{bmatrix} \end{matrix}
 \end{aligned}$$

Example unbalanced (7v1) bi-allelic series

$$\begin{aligned}
 & \text{M.ID} = "0,0,0,0,0,1,0,0" \\
 & \begin{matrix} & k_0 & k_1 \\ \mathbf{M} = & \begin{bmatrix} A & 1 & 0 \\ B & 1 & 0 \\ C & 1 & 0 \\ D & 1 & 0 \\ E & 1 & 0 \\ F & 0 & 1 \\ G & 1 & 0 \\ H & 1 & 0 \end{bmatrix} \end{matrix}
 \end{aligned}$$

$$\begin{aligned}
 & \text{M.ID} = "0,1,0,0,0,0,0,0" \\
 & \begin{matrix} & k_0 & k_1 \\ \mathbf{M} = & \begin{bmatrix} A & 1 & 0 \\ B & 0 & 1 \\ C & 1 & 0 \\ D & 1 & 0 \\ E & 1 & 0 \\ F & 1 & 0 \\ G & 1 & 0 \\ H & 1 & 0 \end{bmatrix} \end{matrix}
 \end{aligned}$$

Example tri-allelic series

$$\begin{aligned}
 & \text{M.ID} = "0,0,1,2,2,0,2,0" \\
 & \begin{matrix} & k_0 & k_1 & k_2 \\ \mathbf{M} = & \begin{bmatrix} A & 1 & 0 & 0 \\ B & 1 & 0 & 0 \\ C & 0 & 1 & 0 \\ D & 0 & 0 & 1 \\ E & 0 & 0 & 1 \\ F & 1 & 0 & 0 \\ G & 0 & 0 & 1 \\ H & 1 & 0 & 0 \end{bmatrix} \end{matrix}
 \end{aligned}$$

$$\begin{aligned}
 & \text{M.ID} = "0,1,0,0,0,0,2,2" \\
 & \begin{matrix} & k_0 & k_1 & k_2 \\ \mathbf{M} = & \begin{bmatrix} A & 1 & 0 & 0 \\ B & 0 & 1 & 0 \\ C & 1 & 0 & 0 \\ D & 1 & 0 & 0 \\ E & 1 & 0 & 0 \\ F & 1 & 0 & 0 \\ G & 0 & 0 & 1 \\ H & 0 & 0 & 1 \end{bmatrix} \end{matrix}
 \end{aligned}$$

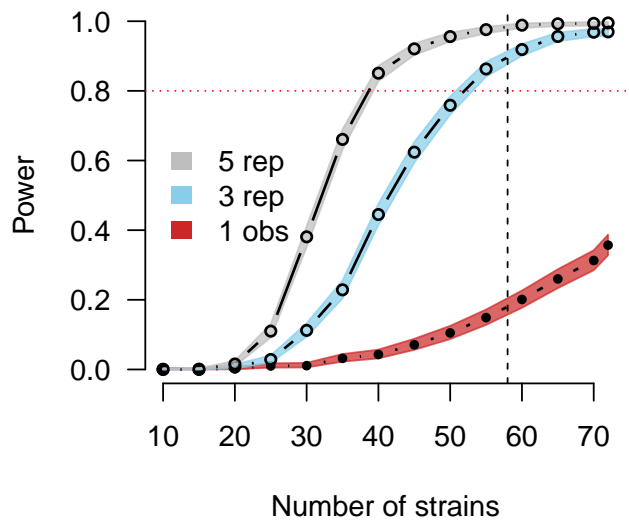


Figure S1 Power estimates for experiments with three and five replicates interpolated from estimates from only a single observation per CC strain. Power curves correspond to a QTL with effect size of 30% and two functional alleles. QTL effect sizes for experiments with replicates are adjusted based on **Eq 4**, allowing for results from single observation simulations to be projected into experiments with replicates. Pre-computed power estimates for single observation simulations are stored in SPARCC and can conveniently be extrapolated into other settings, as is demonstrated here. The horizontal red dotted line marks 80% power. The vertical black dashed line marks 58 strains, which is currently the number of unrelated strains available from UNC. Closed circles represent power estimates that were directly evaluated. Open circles represent power estimates that were interpolated from single observation results.

Supplemental Tables and Figures

Table S1 QTL mapping power in the Collaborative Cross based on QTL effect sizes in the mapping population (Definition DAMB)

QTL			Power								
			30 strains			50 strains			72 strains		
1 obs ^a	3 rep ^b	5 rep ^b	2 alleles	3 alleles	8 alleles	2 alleles	3 alleles	8 alleles	2 alleles	3 alleles	8 alleles
0.01	0.003	0.002	0.000	0.002	0.001	0.000	0.000	0.000	0.000	0.003	0.000
0.05	0.017	0.010	0.000	0.001	0.001	0.000	0.002	0.000	0.005	0.003	0.004
0.1	0.036	0.022	0.004	0.002	0.000	0.009	0.007	0.003	0.015	0.016	0.018
0.15	0.056	0.034	0.002	0.006	0.001	0.016	0.017	0.018	0.056	0.051	0.046
0.2	0.077	0.048	0.008	0.007	0.003	0.032	0.038	0.032	0.119	0.141	0.135
0.25	0.100	0.062	0.009	0.005	0.008	0.071	0.066	0.088	0.264	0.264	0.281
0.3	0.125	0.079	0.015	0.013	0.014	0.141	0.120	0.134	0.460	0.466	0.492
0.35	0.152	0.097	0.028	0.029	0.030	0.234	0.229	0.262	0.695	0.664	0.684
0.4	0.182	0.118	0.045	0.038	0.040	0.415	0.376	0.413	0.854	0.848	0.854
0.45	0.214	0.141	0.082	0.074	0.078	0.603	0.594	0.620	0.958	0.964	0.974
0.5	0.250	0.167	0.136	0.134	0.143	0.769	0.783	0.783	0.996	0.998	0.999
0.55	0.289	0.196	0.198	0.204	0.248	0.911	0.922	0.924	1.000	0.999	1.000
0.6	0.333	0.231	0.334	0.331	0.328	0.985	0.980	0.994	0.999	1.000	1.000
0.65	0.382	0.271	0.519	0.489	0.534	0.998	0.995	0.999	0.999	1.000	1.000
0.7	0.438	0.318	0.707	0.703	0.756	0.998	0.999	1.000	0.999	0.999	1.000
0.75	0.500	0.375	0.866	0.864	0.914	0.998	0.998	1.000	0.999	1.000	1.000
0.8	0.571	0.444	0.940	0.954	0.979	0.996	0.998	1.000	1.000	1.000	1.000
0.85	0.654	0.531	0.962	0.967	0.995	0.998	0.999	1.000	1.000	0.999	1.000
0.9	0.750	0.643	0.978	0.981	0.998	0.999	0.998	1.000	1.000	1.000	1.000
0.95	0.864	0.792	0.970	0.988	0.999	0.999	0.999	1.000	0.999	1.000	1.000

^a Convert QTL effect sizes from experiments with replicates to mean scale with Eq 4.

^b Based on no background strain effect.

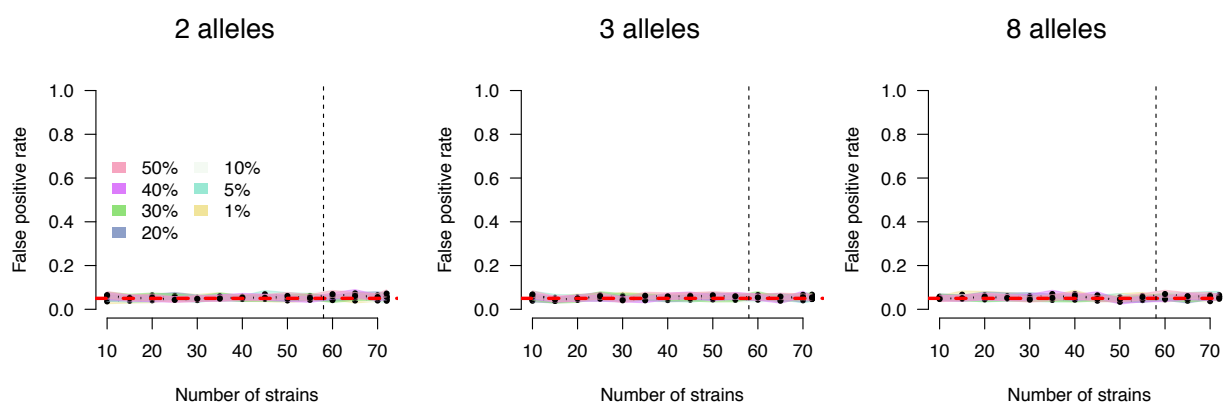
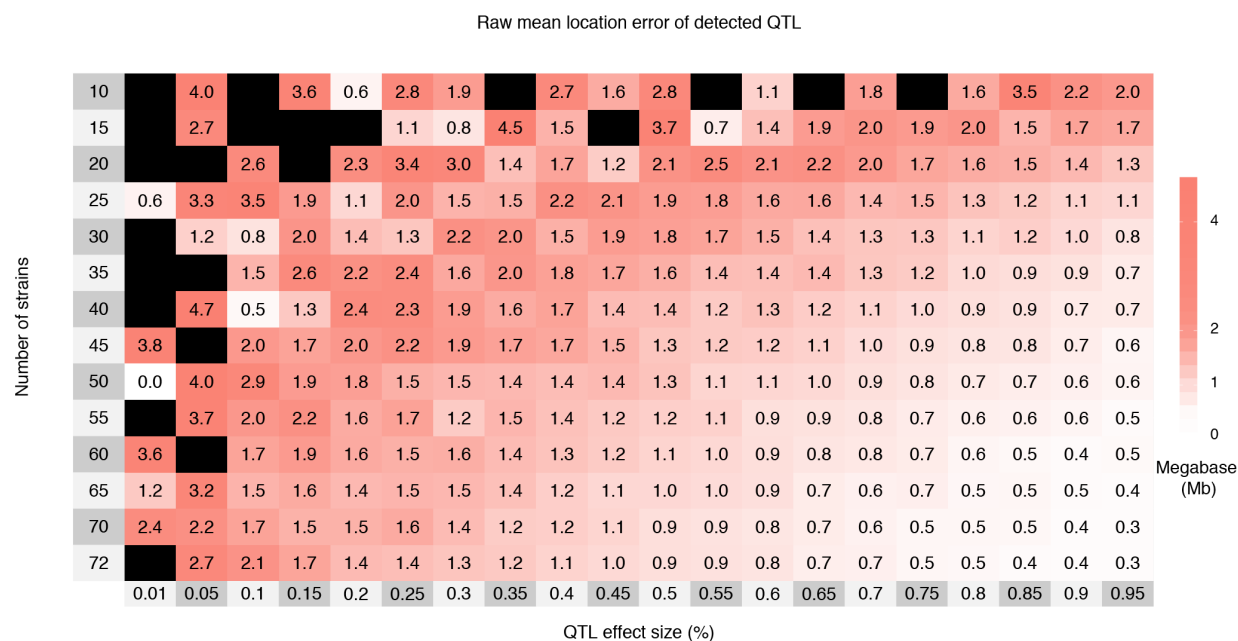


Figure S2 False positive rate (FPR) based on 1,000 simulations per setting with respect to number of CC strains, stratified by the number of functional alleles. The horizontal red dashed line marks the 5% type I error (false positive) rate. CC strains and loci were varied in simulations, resulting in false positive rates that average over loci and strain combinations. Confidence intervals were calculated based on Jeffreys interval (Brown *et al.* 2001) for a binomial proportion. Plots, left to right, correspond to two, three, and eight functional alleles. The FPR represents the probability that any QTL is detected on chromosomes other than the chromosome on which the simulated QTL is located. The significance thresholds maintain the desired type I error rate of 0.05. As expected, the allelic series does not appear to influence FPR. The vertical black dashed line marks 58 strains, which is currently the number of unrelated strains available from UNC.

A



B

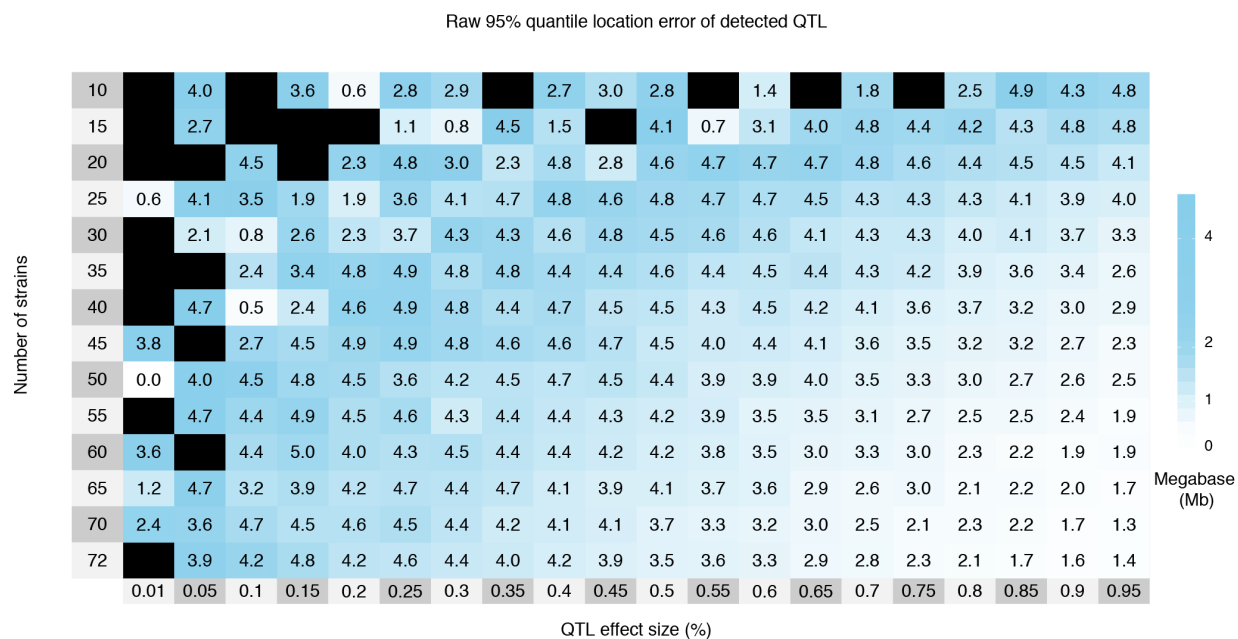


Figure S3 The raw mean (A) and 95% quantile (B) of the location error, the distance in Mb between the detected and simulated QTL, by effect size and number of strains for 1,000 simulations of each setting. These simulations are based on Definition DAMB with an eight allele QTL, and only a single observation per strain. Cells are colored red to white with decreasing mean and blue to white with decreasing 95% quantile. Black cells represent the case in which no simulated QTL were detected. Estimates from poorly-powered settings are more likely to be unobserved or unstable from low detection. Regularized measurements are provided in **Figure 4**. Increasing the number of strains reduces both the mean and 95% quantile location error more so than QTL effect size, also shown in **Figure S5**. The maximum possible location error was 5 Mb due to the 10 Mb window centered around the true QTL position used for detecting QTL.

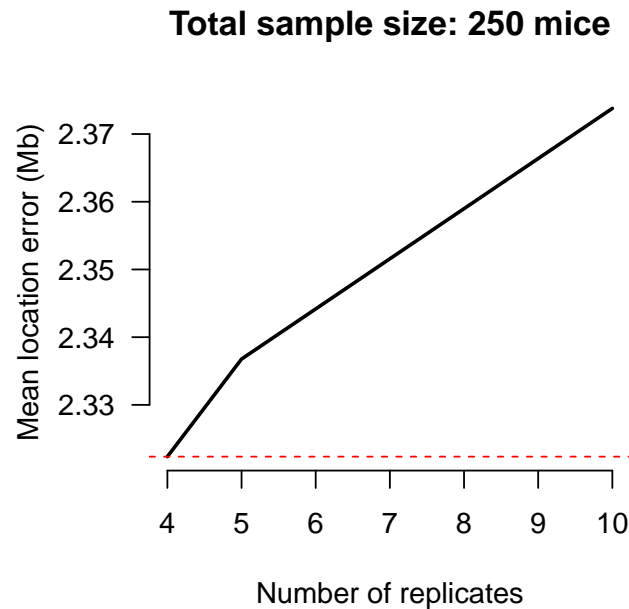


Figure S4 Mean location error of detected QTL increases with the number of replicates while keeping total sample size fixed. Estimates are based on linear interpolation from dense simulations using Definition B with single observations per strains. The total number of mice and the QTL effect size are fixed at 250 and 50%, respectively. The red dotted line highlights that the lowest mean location error occurs at 4, the lowest number of replicates possible for a sample of 250 mice, given the 72 strains used in the simulations.

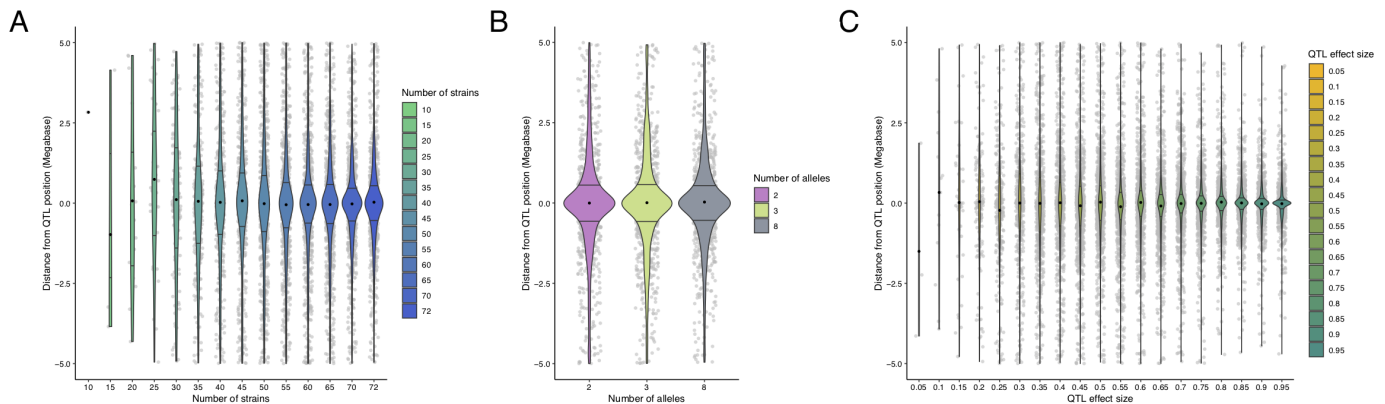


Figure S5 Distributions of the un-regularized location error, by number of strains (A), number of alleles (B), and QTL effect size (C). Observed distances are between -5 and 5 Mb due to the 10 Mb window centered around the simulated QTL that was used for QTL detection in the large scale results. Gray dots represent the distances for a single simulations. The colored violin plots represent the distribution of distances across the simulations. The black dot marks the mean location error for each category. Horizontal lines represent the 25th and 75th quantiles. (A) With QTL effect size fixed at 50% and the number of alleles at 8, as the number of CC strains increases, the distribution of location error becomes more concentrated around zero, meaning the mapping resolution improves with increasing numbers of strains. (B) With the QTL effect size again fixed at 50% and the number of strains fixed at 72, the distribution of distances does not appear to differ based on the number of functional alleles. (C) With the number of strains fixed at 72 and the number of alleles fixed at 8, as the QTL effect size increase, the distribution of distances becomes more concentrated around zero. These simulations are based on Definition B and single observations per strain. See **Figures 4** and **S3** for specific estimates of location error over different settings of QTL effect size and numbers of strains.