

Method of Analysis of Multi-parent Mapping Populations Affects Detection of QTL

Odell, S. G.^{*,1}, Praud, S.[†], Ross-Ibarra, J.^{*,‡} and Runcie, D.^{*,‡}

^{*}Dept. of Plant Sciences and Center for Population Biology, University of California, Davis, CA, USA, [†]Limagrain, Chappes, France, [‡]Genome Center, University of California, Davis, CA, USA

ABSTRACT The search for quantitative trait loci (QTL) that explain complex traits such as yield and drought tolerance has been ongoing in all crops. Methods such as bi-parental QTL mapping and genome-wide association studies (GWAS) each have their own advantages and limitations. Multi-parent advanced generation inter-crossing (MAGIC) contain more recombination events and genetic diversity than bi-parental mapping populations and reduce the confounding effect of population structure that is an issue in association mapping populations. Here we discuss the results of using a MAGIC population of doubled haploid (DH) maize lines created from 16 diverse founders to perform QTL mapping, comparing QTL identified using a 600K SNP array to those found using founder probabilities and haplotype probabilities generated by determining the regions of the MAGIC DH lines that were derived from the 16 founders and by identifying regions of identity-by-descent (IBD) between the 16 founders, respectively. The three methods have differing power and resolution for detecting QTL for a variety of agronomic traits. This highlights the importance of considering different approaches to analyzing genotypic datasets, and shows the limitations of binary SNP data for identifying multi-allelic QTL.

KEYWORDS QTL, MAGIC

Introduction

The study of evolutionary quantitative genetics requires the ability to link differences in phenotype to genotypic variation. Natural and artificial selection acts on phenotypes, but only heritable phenotypic variation will result in changes in population means. Maize presents an excellent model organism to answer these questions due to the combination of extensive genetic and phenotypic resources, and the ability to create mapping populations. In addition, maize is one of the most widely produced crops in the world and is a major source of calories for millions of people. Decades of research into maize genetics have resulted in the identification of many quantitative trait loci (QTL) that explain variation in phenotypes such as yield, flowering time, and plant height ((Buckler *et al.* 2009);(Beavis *et al.* 1991);(Martinez *et al.* 2016)). Such traits are extremely agronomically important. They are also crucial plant phenotypes in terms of fitness and local adaptation. Researchers have been able to discover large-effect QTL for a number of agronomic traits through the use of different types of mapping populations (Wallace *et al.* 2014). The choice of any population comes with associated advantages and limitations. In particular, they tend to vary in two main characteristics: (1) their ability to capture genetic diversity and (2) their power to detect QTL of small effect. Multi-parent Advanced Generation Intercross (MAGIC) populations have been used in breeding to increase the genetic diversity and number of recombination events included in a mapping population compared to biparental populations (Huang *et al.* 2012)[give more exam-

ple citations]. Compared to genome-wide association panels, MAGIC populations have more power to detect rare alleles (i.e. alleles that are only present in one of the parents) and can better estimate allelic effects because the crossing scheme increases the frequency of all parental alleles to be approximately equal. Simulations of an 8-parent MAGIC population showed that sample sizes of 300 could detect QTL accounting for 12% of variance with a power of 82% (Dell'Acqua *et al.* 2015). Lastly, a MAGIC population avoids confounding due to population structure that is encountered with GWAS because the pedigree of the lines is known. In this study, we utilized a MAGIC population of about 400 doubled-haploid lines derived from 16 inbred maize parents developed by Limagrain to understand how methods of representing genotypic data can impact the identification of QTL. Extensive genetic resources already exist for maize, but do not possess the same diversity and statistical power as the Biogemma MAGIC population. A maize nested association mapping (NAM) population exists, consisting of RIL populations derived from 25 inbred parents crossed to B73 (Yu *et al.* 2008). Only two inbred parents overlap between the NAM and Biogemma MAGIC populations (B73 and Oh43), and compared to the NAM, the MAGIC population can have similar power to the NAM using half the number of samples (Dell'Acqua *et al.* 2015). Likewise, another maize MAGIC population has previously been created, which overlaps by three parents (A632, B73, and B96) (Dell'Acqua *et al.* 2015). However, the previous MAGIC is derived from 8 inbred maize parents instead of 16, and consists of RILs, not doubled haploids, so some residual heterozygosity may exist. For these reasons, the Biogemma MAGIC population has great potential to reveal new insights into the genetic control of quantitative traits in maize. It also serves a re-

¹ Dept. of Plant Sciences, University of California, Davis, CA, USA E-mail: sgodell@ucdavis.edu

liable standard for QTL mapping method comparison because DH lines were used and the crossing scheme ensured that there is a reasonably even distribution of the 16 founders within the population (supplemental figure?)

A good introduction is very important.

Materials and Methods

Genotype Data

The MAGIC population was derived from 16 inbred maize parents representing the diversity of European flint and U.S. dent heterotic groups. The 16 founder lines were crossed in a funnel crossing scheme, and then the resulting synthetic population was intercrossed for 3 generations with around 2000 individuals per cycle (Figure 1A). Finally, 800 lines were selected from the synthetic population to create doubled haploids (DH), resulting in 550 MAGIC DH lines at the end of the process. The 16 founder lines and the MAGIC DH lines were all genotyped with the Affymetrix 600K Axiom SNP array. In addition, the 16 founder lines were sequenced with Illumina short-read sequencing to a depth of [?]x, resulting in 45.4 millions SNPs and 5.4 indels after filtering using GATK best practices.

Phenotype Data

The MAGIC DH lines were crossed to a tester MBS84 to produce 344 hybrids. Due to variation in flowering time, a subset of the lines could not be crossed to the tester. The hybrids were grown in the field in Blois, France in 2014. For each genotype, two blocks of around 80 plants were grown under well-watered conditions. Measured phenotypes included days to anthesis (DTA), days to silking (DTS), plant height, percent harvest grain moisture, grain yield, and thousand kernel weight (adjusted to 15% humidity), where values were averaged over blocks. Both flowering time phenotypes were measured as the sum of degree days since sowing with a base temperature of 6°C (48°F). Days to anthesis was considered as the growing degree days until 50% of plants in a block were flowering at 25% of the central tassel spike.

Calculation and Validation of Founder Probabilities

The package R/qtl2 (Broman *et al.* 2019) was used to determine founder probabilities of the DH lines using the 600K genotype data and the cross type “riself16”. Due to the fact that the actual crossing scheme and the cross type input into R/qtl2 differed, we wanted to assess the accuracy of the founder probabilities. This was done by simulating lines using the actual crossing scheme and assessing the performance of the `calc_genoprob` function of R/qtl2 in correctly identifying the founder genotype (Figure 1). We developed an R package, `magicsim` (<https://github.com/sarahodell/magicsim>), to simulate the lines using the maize genetic map from Ogut *et al.* 2015 to generate approximate recombination rates across the chromosome. For 400 simulated lines, 99.6% of SNPs were correctly assigned to the founder genotype.

Calculation of Haplotype Probabilities

The use of founder probabilities makes the assumption that all 16 founders have distinct haplotypes. This assumption is not realistic, especially considering the known varying degrees of relatedness between the founders. The identification of regions of shared genetic sequence between founder pairs would allow for the collapsing of founder probabilities into haplotype

probabilities. These haplotype probabilities have the potential to increase statistical power by reducing the number of tests performed in QTL mapping. Areas of uncertainty in founder probabilities of the DH lines were associated with regions of identity-by-descent (IBD) between two or more founder lines in that region of the chromosome. Identity-by-state (IBS) was measured from the 600K genotype data of the founders using the software GERMLINE with the parameters “-wextend”, “-min_m 1”, and “-err_hom 4” as a proxy for IBD. In order to assess the quality of predicted IBD regions, we also identified IBD regions using the founder WGS data. This was done using two separate softwares. GERMLINE was run using the same parameters as above on the WGS data with all sites with missing data removed, resulting in a total of 16.9 million SNPs for comparison of sequence similarity. The software IBDSeq was run in addition using the full 45.4 million sites in the WGS data. Comparison of overlap between pairwise IBD regions among the three methods showed [... significant overlap between the 600K and WGS data... actual numbers]. Haplotype blocks were created by grouping regions of pairwise IBD across the 16 founders. Within blocks, the founder probabilities for founders that were in IBD were summed to obtain haplotype probabilities. This resulted in haplotype blocks with the number of unique haplotypes within blocks ranging from [2?] to 16.

Association and QTL Mapping

The R package GridLMM [] was used to run association mapping using the three different methods of representing the genotype data. The function `GridLMM_ML` was used with the “ML” option. The model for each was

$$Y = X\beta + Zu + \epsilon \quad (1)$$

where Y is the response variable, X is the genotype matrix, β , is the effect size, Z is the design matrix, u is the random effects, and ϵ is the error. Significance cutoffs for p-values were obtained using permutation testing, taking the 5% cutoff from 1000 randomized permutations for each method. Code for the analyses can be found [github link].

Method Comparison

The results of the three methods of identifying QTL were compared using two main criteria: (i) presence or absence of identified QTL peaks and (ii) the resolution of those QTL peaks. The GWAS methods should be most powerful at identifying QTL for which the causal variant is biallelic and the tagged SNP is in tight LD with the causal variant. However, for multi-allelic QTL or QTL for which LD is low between tagged SNPs, this method begins to lose power. Founder probabilities increase the odds of detecting both QTL that are multi-allelic and QTL whose causal variant is not in tight LD with any one tagged SNP [Supplemental Figure?]. Lastly, haplotype probabilities potentially improve on the power of founder probabilities to detect QTL that meet the above criteria by reducing the number of test. However, haplotype probabilities also may obscure the signal of some QTL if founders are called as being in IBD with one another when they actually differ for the causal variant. Due to the fact that the founder and haplotype probabilities take into account recent recombination events, whereas the GWAS method only uses historical recombination, we predicted that founder and haplotype mapping would result in higher resolution around QTL peaks. Higher resolution QTL are ideal in that it makes it easier to narrow down candidate genes and potential causal variants when the significant window is smaller.

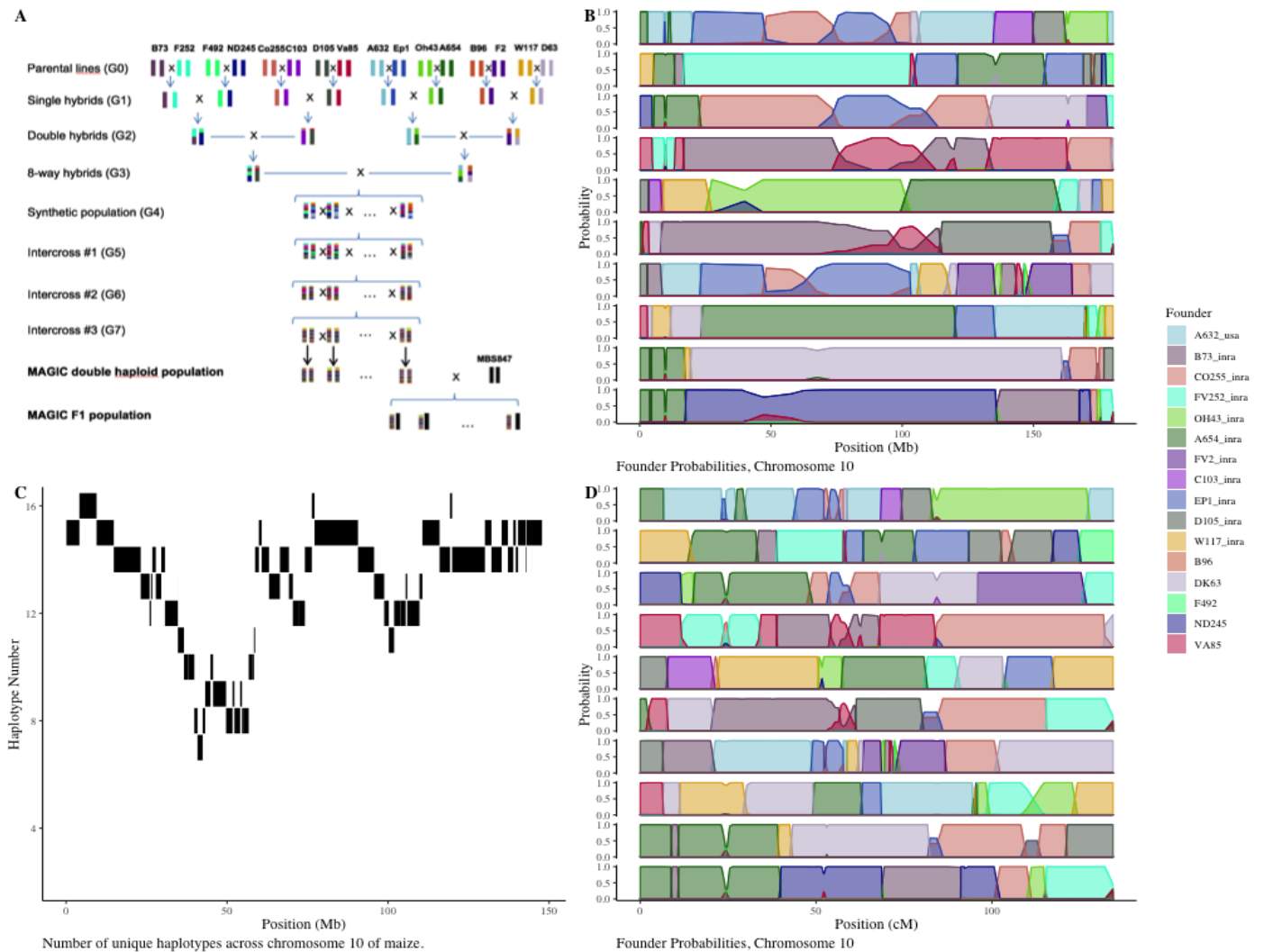


Figure 1 Structure, diversity, and founder representation of the MAGIC population. (A) The crossing scheme of the MAGIC population. (B) Founder probabilities for 10 MAGIC DH lines on chromosome 10 in physical distance. (C) The number of unique haplotype groups among the 16 founder lines across chromosome 10. (D) Founder probabilities for 10 MAGIC DH lines on chromosome 10 in genetic distance.

Results

MAGIC Population

Analysis of the MAGIC population showed that the representation of the 16 founders in the MAGIC DH lines was relatively even, with the highest percentage founder, A654, representing 6.741% and the lowest percentage founder, EP1, representing 5.237% (Supplemental Figure). Across all regions of the chromosomes, for 99% (?) of markers in the 600K array, at least one DH line was derived from each of the 16 founders with high confidence (greater than 0.9 probability) (Supplemental Figure). One notable exception was on chromosome 6, where a region [location] had no representation from [founders]. This may be due to these founders being in close IBD with another founder, resulting in uncertainty in the founder probabilities. Another possible explanation is that there was selection against these founders during the breeding process. This region approximately corresponds to the NOR region [citation].

The founder probabilities determined using R/qtl2 were able to assign founders to the DH lines with high confidence (>0.90)

for [percentage] of the 10 chromosomes of maize.

Simulation and Validation of Founder Probabilities

Using the package we created, magicsim, we simulated chromosome 10 for 1000 MAGIC DH lines following the crossing scheme used for the actual population. We then ran R/qtl2 on these simulated lines to obtain founder probabilities and compared these probabilities to the known founder identities of the simulated lines. In total, for sites for which R/qtl2 was able to assign founder probabilities, the maximum probability founder matched the actual founder for 99.8% of sites. This reinforced our confidence in the founder probabilities obtained from the actual data.

Phenotype Data

The distribution of the six measured phenotypes were approximately normal (Figure 2).

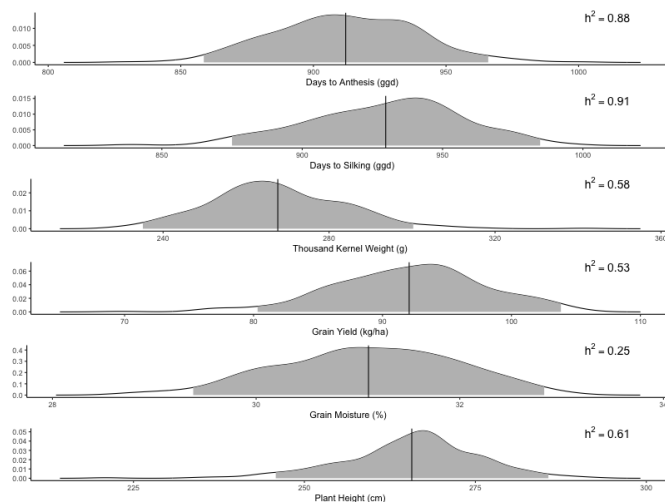


Figure 2 Distribution of phenotypes The density plots of the six measured phenotypes. The vertical bar represents the mean, and the grey shading shows two standard deviations. The heritability of each trait is shown on the right.

QTL Mapping and Association Mapping

The three methods varied both in their ability to identify QTL and the resolution of those QTL.

Variation around *vg1*

A previously characterized QTL, *vg1* is associated with variation in flowering time [citations]. An earlier flowering phenotype is strongly correlated with a MITE insertion about 70kb upstream of the flowering time regulator, *ZmRAP2.7*, an *APETALA*-like transcription factor [].

Discussion

Some more test

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

another figure

This is a full size figure (Figure S1) just add stars to the figure

Acknowledgments

We acknowledge the support of our coffee maker that made this work possible

References

- Beavis, W. D., D. Grant, M. Albertsen, and R. Fincher, 1991 Quantitative trait loci for plant height in four maize populations and their associations with qualitative genetic loci. *Theoretical and Applied Genetics* **83**: 141–145.
- Broman, K. W., D. M. Gatti, P. Simecek, N. A. Furlotte, P. Prins, *et al.*, 2019 R/qt2: Software for mapping quantitative trait loci with high-dimensional data and multiparent populations. *Genetics* **211**: 495–502.
- Buckler, E. S., J. B. Holland, P. J. Bradbury, C. B. Acharya, P. J. Brown, *et al.*, 2009 The genetic architecture of maize flowering time. *Science* **325**: 714–718.
- 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.
- Huang, B. E., A. W. George, K. L. Forrest, A. Kilian, M. J. Hayden, *et al.*, 2012 A multiparent advanced generation inter-cross population for genetic analysis in wheat. *Plant Biotechnology Journal* **10**: 826–839.
- Martinez, A. K., J. M. Soriano, R. Tuberosa, R. Koumproglou, T. Jahrmann, *et al.*, 2016 Yield qtlome distribution correlates with gene density in maize. *Plant Science* **242**: 300–309.
- Wallace, J. G., P. J. Bradbury, N. Zhang, Y. Gibon, M. Stitt, *et al.*, 2014 Association mapping across numerous traits reveals patterns of functional variation in maize. *PLOS Genetics* **10**: e1004845.
- 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.

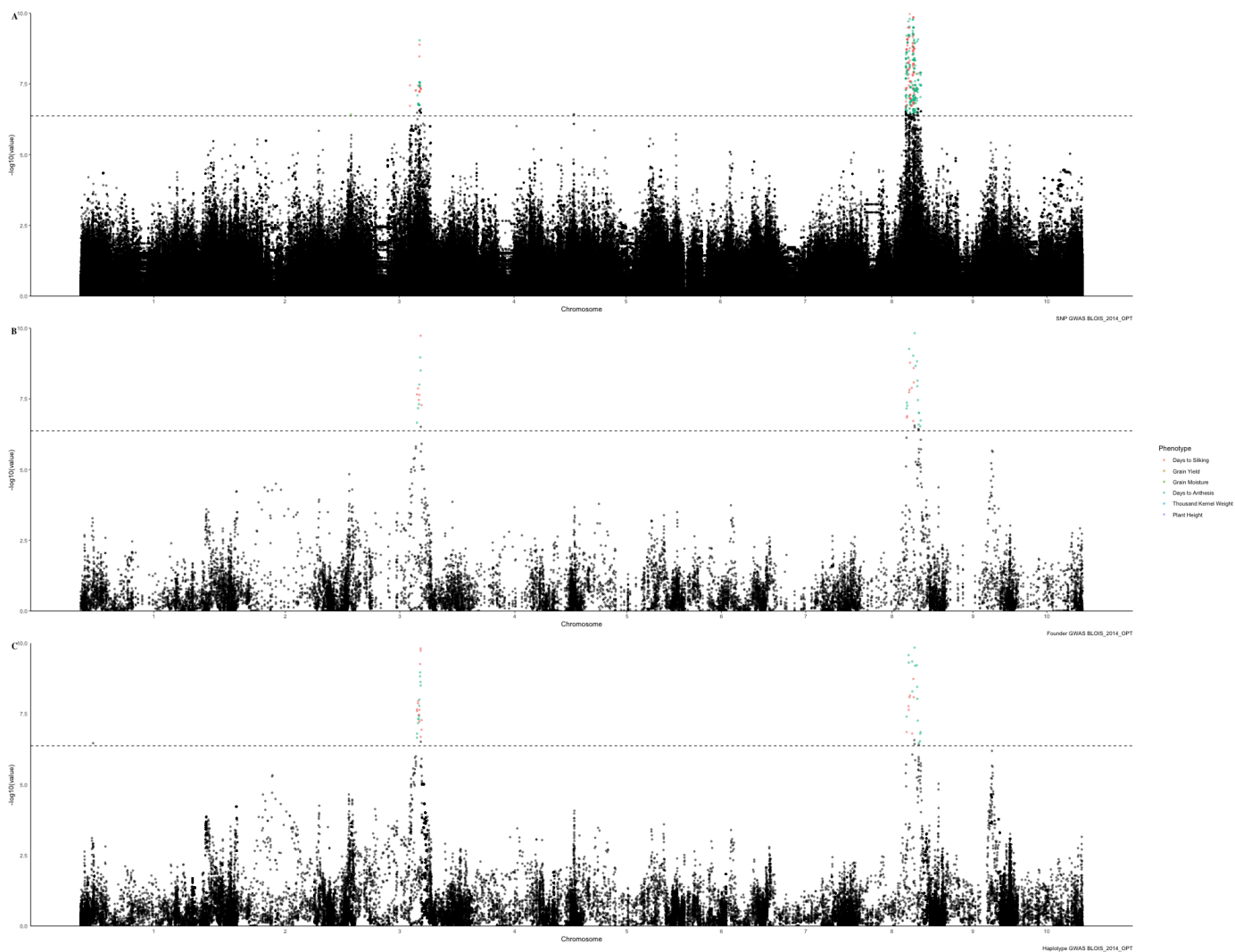


Figure 3 Results of three methods of QTL identification Colored points represent significant SNPs above the 5% significance threshold from 1000 random permutations **top** GWAS results using the 600K SNP array. **middle** Results from QTL mapping using founder probabilities **bottom** Results from QTL mapping using haplotype probabilities

Supplement

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

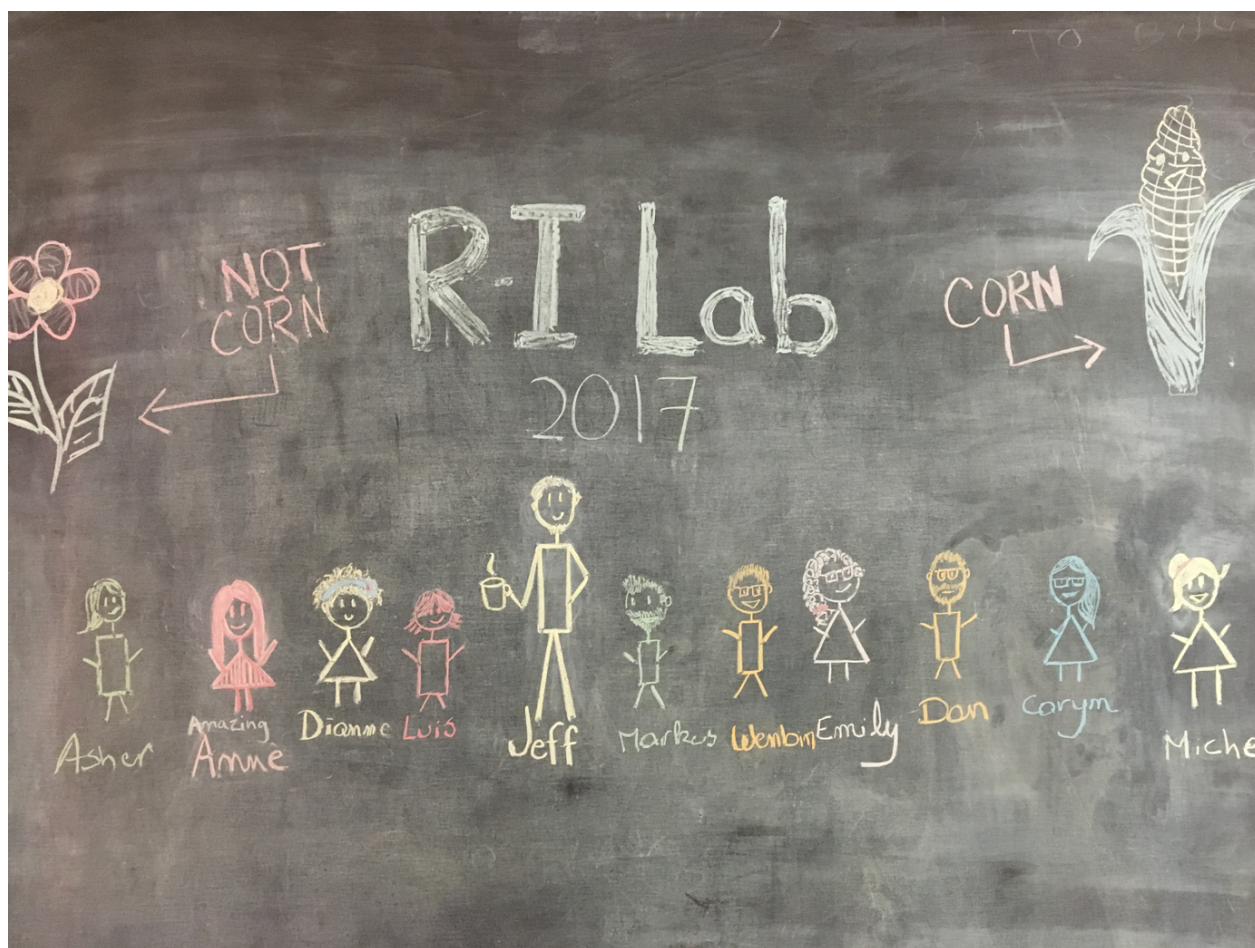


Figure S1 Supplemental figure Test test test

Table S1 Shrink a large table to fit the page

Parameter	Description
Adaptation	Trait related parameters
Time to optimum	Generations until new optimum is reached
Adaptation rate (haldane)	Adaptation rate until new optimum is reached. Calculated as $rate(h) = \frac{\frac{\ln(x_2)}{sd_{x12}} - \frac{\ln(x_1)}{sd_{x12}}}{t_2 - t_1}$
Final genetic variance	Genetic variance in the final generation
Fixations	Mutations that fix after the optimum shift
From new mutations (#)	Sum of fixed mutations in the final population that were already segregating before the optimum shift
From standing variation (#)	Sum of fixed mutations in the final population that arose after the optimum shift
Max. effect size	Maximal effect size of all fixations
Mean effect size	Mean effect size of all fixations
Mean effect size of negative fixations	Mean effect size of negative mutations
Mean effect size of positive fixations	Mean effect size of positive mutations
Mean emergence time	Mean generation when a mutation arose that fixed in the last 0.1 N generations
Mean fixation time	Mean generation in which a mutation fixed
Min. effect size	Minimal effect size of all fixations
Negative (#)	Sum of fixed mutations with negative effects in the final population
New/standing fixations	Ratio of mutations from new mutations vs. standing mutations
Proportion negative	Proportion of negative fixations from all mutations
Positive (#)	Sum of fixed mutations with positive effects in the final population
SD of effect sizes	Standard deviation of effect sizes of all fixations
SD of negative effect sizes	Standard deviation of effect sizes of negative fixations
SD of positive effect sizes	Standard deviation of effect sizes of positive fixations
Total (#)	Sum of fixed mutations in the final population
Sweeps	Mutations that fix faster than 99% of neutral fixations
Hard sweeps (#)	Sum of selective sweeps from new mutations
Proportion of hard sweeps	Porportion of hard selective sweeps of all selective sweeps
Proportion of sweeps from standing	Proportion of selective sweeps from ststanding variation of all selection sweeps
Sweeps (#)	Sum of selective sweeps
Sweeps from standing variation (#)	Sum of selective sweeps from mutations that were already segregating before the optimum shift
Sweeps/fixations	Ratio of sweeps vs. fixations
Segregating sites	Mutations that segregate in the final generation
Max. effect size	Maximal effect size of segregating sites
Mean effect size	Mean effect size of segregating sites
Mean effect size of negative sites	Mean effect size of segregating sites with negative effects
Mean effect size of positive sites	Mean effect size of segregating sites with positive effects
Mean frequency of all sites	Mean allele frequency of segregating sites
Mean frequency of negative sites	Mean allele frequency of segregating sites with negative effects
Mean frequency of positive sites	Mean allele frequency of segregating sites with positive effects
Min. effect size	Minimal effect size of segregating sites
Negative (#)	Sum of segregating sites with negative effect
Positive (#)	Sum of segregating sites with positive effect
Proportion of negative sites	Proportion of segregating sites with negative effect of all segregating sites
Standard deviation of effect sizes	Standard deviation of effect sizes of all segregating sites
Total (#)	Sum segregating sites in the final generation