# UC Riverside
## UC Riverside Previously Published Works

**Title**

Genetic architecture of nonadditive inheritance in Arabidopsis thaliana hybrids

**Permalink**

https://escholarship.org/uc/item/4xr1p4hk

**Journal**

Proceedings of the National Academy of Sciences of the United States of America, 113(46)

**ISSN**

0027-8424

**Authors**

Seymour, Danelle K
Chae, Eunyoung
Grimm, Dominik G
et al.

**Publication Date**

2016-11-15

**DOI**

10.1073/pnas.1615268113

Peer reviewed

# Genetic architecture of nonadditive inheritance in *Arabidopsis thaliana* hybrids

Danelle K. Seymour[a,1], Eunyoung Chae[a,1], Dominik G. Grimm[b,c,1], Carmen Martín Pizarro[a,2], Anette Habring-Müller[a], François Vasseur[a], Barbara Rakitsch[b,d], Karsten M. Borgwardt[b,c], Daniel Koenig[a,3], and Detlef Weigel[a,4]

[a]Department of Molecular Biology, Max Planck Institute for Developmental Biology, 72076 Tubingen, Germany; [b]Machine Learning and Computational Biology Research Group, Max Planck Institute for Developmental Biology, Max Planck Institute for Intelligent Systems, 72076 Tubingen, Germany; [c]Machine Learning and Computational Biology Laboratory, Department of Biosystems Science and Engineering, Eidgenössische Technische Hochschule Zürich, 4058 Basel, Switzerland; and [d]European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Trust Genome Campus, Cambridge, CB10 1SD, United Kingdom

The ubiquity of nonparental hybrid phenotypes, such as hybrid vigor and hybrid inferiority, has interested biologists for over a century and is of considerable agricultural importance. Although examples of both phenomena have been subject to intense investigation, no general model for the molecular basis of nonadditive genetic variance has emerged, and prediction of hybrid phenotypes from parental information continues to be a challenge. Here we explore the genetics of hybrid phenotype in 435 *Arabidopsis thaliana* individuals derived from intercrosses of 30 parents in a half diallel mating scheme. We find that nonadditive genetic effects are a major component of genetic variation in this population and that the genetic basis of hybrid phenotype can be mapped using genome-wide association (GWA) techniques. Significant loci together can explain as much as 20% of phenotypic variation in the surveyed population and include examples that have both classical dominant and overdominant effects. One candidate region inherited dominantly in the half diallel contains the gene for the MADS-box transcription factor *AGAMOUS-LIKE 50* (*AGL50*), which we show directly to alter flowering time in the predicted manner. Our study not only illustrates the promise of GWA approaches to dissect the genetic architecture underpinning hybrid performance but also demonstrates the contribution of classical dominance to genetic variance.

heterosis | half diallel | GWAS | *Arabidopsis thaliana*

The often observed phenotypic superiority of progeny relative to their parents, or heterosis, is a universal phenomenon and of great importance to plant agriculture. The earliest description of heterosis (also known as hybrid vigor or superiority) dates to Darwin's studies of cross-fertilization in plants. He noticed that intercrossing distantly related individuals gave rise to larger, more vigorous progeny (1). Four decades later, George Shull coined the term "heterosis" (2) for this phenomenon, which he and Edward East had independently described for hybrids of inbred maize in 1908 (3, 4). Heterosis has long been of interest to evolutionary biologists as a potential explanation for the ubiquity of cross-fertilization in plants and animals, but it is also a central component of agricultural breeding programs. The combination of hybrid seed technology and inbred line improvement has driven an unprecedented improvement in maize yield over the past century (5). Despite the economic importance of heterosis and its intensive investigation in a wide spectrum of species, prediction of hybrid performance from parental information remains a major challenge (6).

In the terms of quantitative genetics, hybrid vigor (and its opposite, inferiority) describes a deviation of progeny from the phenotypic mean of the parents. This means that heterosis cannot be explained by the addition of the effects of contributing alleles (7). Nonadditive genetic variance can result from a nonlinear phenotypic effect of alleles at one locus, as in the case of dominant/recessive allele pairs in classical genetics, or from epistatic interactions between loci (reviewed in ref. 8). Three nonmutually exclusive types of intralocus interactions are commonly invoked to

explain heterosis. The overdominance hypothesis postulates that a single mutation in the heterozygous state is causal (3, 4, 9), and it accounts for at least a few cases of heterosis (10–13) and hybrid inferiority (14, 15). The dominance hypothesis suggests that genome-wide complementation of many small-effect, weakly deleterious loci drives hybrid superiority (16–18), but the small effect size of individual loci would make it difficult, if not impossible, to generate direct support for this hypothesis. Finally, the pseudo-verdominance hypothesis also explains heterosis with the complementation of recessive alleles but proposes that when linked in repulsion, such alleles appear overdominant. Outside of these classical hypotheses, some cases of hybrid inferiority (19–22) and hybrid superiority (23–30) have been linked to epistatic interactions between parental alleles.

The availability of completely homozygous natural accessions has made the model plant *Arabidopsis thaliana* an excellent subject for studies of natural variation. Collections of sequenced accessions enable replicated genome-wide association (GWA)

## Significance

Hybrid progeny of inbred parents are often more fit than their parents. Such hybrid vigor, or heterosis, is the focus of many plant breeding programs, and the rewards are evident. Hybrid maize has for many decades accounted for the majority of seed planted each year in North America and Europe. Despite the prevalence of this phenomenon and its agricultural importance, the genetic basis of heterotic traits is still unclear. We have used a large collection of first-generation hybrids in *Arabidopsis thaliana* to characterize the genetics of heterosis in this model plant. We have identified loci that contribute substantially to hybrid vigor and show that a subset of these exhibits classical dominance, an important finding with direct implications for crop improvement.
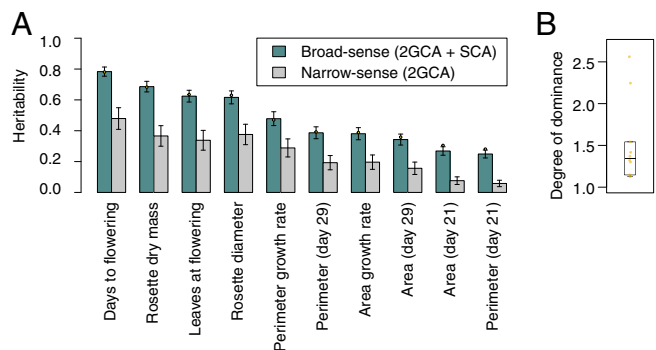
PLANT BIOLOGY

mapping studies across varied environmental conditions (31–33). The drawback of highly inbred lines is that the contribution of dominance to phenotype cannot be assessed directly. However, *Arabidopsis thaliana* outcrossing rates of over 10% have been reported in the field, suggesting that dominance may contribute to phenotypic variation in natural populations (34). Here we explore the magnitude of nonadditive genetic variation in *A. thaliana* using a half diallel intercrossing scheme. This approach was chosen because of its power to separate a line's breeding value (additive contribution) from its performance in a specific cross (nonadditive contribution) (35). Whole-genome resequencing information is available for all 30 parental accessions in our scheme (36), enabling construction of hybrid genotypes and GWA mapping of hybrid phenotypes. We show that nonadditive inheritance is pervasive in *A. thaliana* hybrids, that the genetic basis of such traits can be uncovered using a modified GWA approach, and that a candidate gene underlying an associated peak is sufficient to alter a flowering time trait.

## Results

**Experimental Design and Phenotypic Analyses.** A half diallel was constructed by intercrossing 30 natural accessions of *A. thaliana* (*SI Appendix*, Table S1). These accessions were chosen because they span much of the genetic diversity in the native range of the species, and their genomes have been sequenced (*SI Appendix*, Fig. S1) (36). To facilitate the large number of intercrosses, male sterile lines were generated by artificial miRNA knockdown of the homeotic gene *AP3*, removing the need for manual emasculation (22). Because manual crossing is known to influence trait values of the progeny even when using genetically identical parents (37), we manually self-crossed each parental line using *AP3* knockdown females and wild-type males as controls. This crossing scheme resulted in 435 hybrid genotypes and $2 \times 30$ parental genotypes (both normally and manually selfed lines). These were grown in 16 °C long days in a completely randomized design with five replicates per genotype. Plants were phenotyped for traits related to flowering time [days to flowering (DTF) and leaves on the main shoot at flowering (LTF)] and final rosette size (rosette diameter and rosette dry mass). Additionally, images were taken of young rosettes (21 and 29 days after sowing), and several rosette traits were extracted from these images.

We often observed differences between progeny from natural self-fertilization and progeny from manual fertilization of *AP3* knockdown females with pollen from isogenic siblings (*SI Appendix*, Figs. S2 and S3). Although such differences between otherwise genetically identical individuals have been reported before (37), our much larger dataset demonstrates that the effect is not directional, with progeny of the manual crosses not always being larger than their self-fertilized siblings (*SI Appendix*, Figs. S2 and S3). The artificial miRNA itself is not the source of these differences, because the presence of the transgene explains very little, if any, of the total phenotypic variance (*Materials and Methods*). Instead, discrepancies between these two groups of parental genotypes likely result from strong maternal effects; for example, knockdown of *AP3* in the female parents greatly diminishes fruit production, potentially altering resource allocation. Regardless of the mechanism, the crossing process clearly influenced the phenotypes of resulting progeny. With this in mind, we only used phenotypes from manually crossed parents in our analyses below.

To understand the genetic independence of the measured traits, we estimated their genetic correlation (*SI Appendix*, Fig. S4). Several traits (DTF, LTF, and dry mass) were correlated and thus shared a genetic basis (*SI Appendix*, Fig. S4). The remaining rosette traits were also correlated but were not, or only very weakly, correlated to the flowering time traits, suggesting that the genetic basis of rosette size over time and onset of flowering are largely independent (*SI Appendix*, Fig. S4).



**Fig. 1.** Broad- and narrow-sense heritability estimates. (*A*) Broad- and narrow-sense heritabilities for each trait. Heritabilities were calculated from estimates of GCA and SCA, which were derived from an LMM implemented in SAS. Yellow circle shows the broad-sense heritability reestimated using a linear model in R. (*B*) Boxplot of the degree of dominance for all traits calculated as $\sqrt{\frac{2\sigma_D^2}{\sigma_A^2}}$.

We next sought to estimate the relative contributions of additive and dominance components to overall phenotypic variation in our sample. With diallel designs, one can evaluate the breeding value of each parent, or its general combining ability (GCA). One can also estimate the specific combining ability (SCA) (35). The SCA is a measure of deviation from the breeding values, or the expected performance of a line in a particular hybrid. Because additive and dominance genetic variance can be derived from estimates of GCA and SCA, respectively, it is possible to calculate both narrow- and broad-sense heritability using these designs (*Materials and Methods*). We estimated GCA, SCA, and heritabilities for each trait (Fig. 1*A* and *SI Appendix*, Table S2) using a linear mixed model (LMM) (*Materials and Methods*). Total genetic variance (broad-sense heritability) ranged from 24 to 78% of the total phenotypic variance (Fig. 1*A*). Rosette traits of younger plants estimated from images seemed to have much lower broad-sense heritability than adult traits. Despite the large range of broad-sense heritability estimates, nonadditive inheritance contributed substantially to the observed hybrid traits. The degree of dominance, calculated as a ratio of nonadditive to additive genetic variation, was greater than 1 for all traits, and such values can only be explained by either overdominance or pseudooverdominance (Fig. 1*B*) (38). The significant contribution of nonadditivity to overall genetic variation suggests that our experimental system is ideal for understanding the factors underlying nonadditive inheritance.

**Model Selection, Simulation of Phenotypes, and Power Analyses.** Our next goal was to identify loci that were contributing to dominance and heterosis using GWA. Typically, GWA studies of continuous traits search for a linear relationship between genotypic class and a trait of interest. For binary traits, such as those frequently used in human disease case-control studies, more complex genetic models, including dominance and overdominance, can be explicitly tested (39, 40). Although this is rarely done in the analysis of continuous traits, a few studies have reported associations of heterotic traits with heterozygosity (41–45). We selected two linear mixed models to search for associations between genotype and phenotype using FaSTLMM software in the easyGWAS framework (46, 47). The first model used a standard linear additive SNP encoding, where the homozygous major allele was represented as 0, the heterozygous as 1, and the homozygous minor allele as 2; we refer to it as the "additive model." The second model, referred to as the "overdominant model," used a nonstandard SNP encoding, where both homozygous classes were represented as 0 and the heterozygous genotype as 1. The genetic similarity between individuals was estimated by computing

the realized relationship kinship matrix using SNP encodings specific for each model (48).

In addition to fitting two different models to our data, we chose to search for association of variants not only with estimated trait means, but we also estimated the discrepancy of an observed hybrid phenotype from its midparent performance [midparent heterosis (MPH)] for each hybrid–parent trait combination (7). Because of the bidirectional discrepancy between self-fertilized and manually fertilized parental genotypes, phenotypic means from parental genotypes produced by manual crosses were used to estimate MPH (*SI Appendix*, Figs. S2, S3, and S5). By mapping MPH, we were able to increase the sensitivity to detect nonadditive loci. The three tested model–phenotypic component combinations are summarized in *SI Appendix*, Fig. S6.
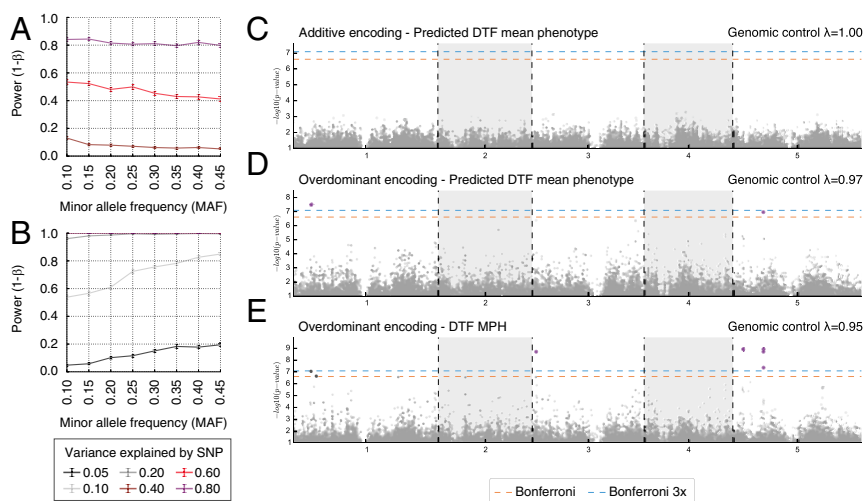
Because our design was different from those used in previous *A. thaliana* GWA studies, which used only homozygous genotypes, we used simulations to estimate our power to detect loci with additive or dominance effects. It turned out that the additive model was extremely underpowered in this dataset regardless of the variance explained, which is a proxy for effect size of individual SNPs, or of allele frequency of the causal SNP (Fig. 2*A*). This could be the result of the correlation of such sites with population structure. Alternatively, it may be caused by fewer available degrees of freedom, due to the limited sample size of the source population ($n = 28$). Simulations also showed that in contrast to the additive model, the overdominant model had sufficient power to detect associations with SNPs that explained a range of variances and that had different minor allele frequencies (Fig. 2*B*), emphasizing the importance of the diallel design in our study.

**GWA Mapping of Additive and Nonadditive Inheritance.** In silico $F_1$ genotypes were constructed by combining known parental genotypes (36). Informative sites were required to have complete information, with a minor allele frequency of at least 10% in the diallel. Due to the limited genetic diversity in the founding parents, many positions were in complete linkage disequilibrium (LD) across chromosomes. These as well as positions in LD with

10 or more other sites were excluded from GWA tests, leaving 204,753 sites segregating in the diallel population (*SI Appendix*, Figs. S7 and S8).

We used the three approaches described above to identify informative SNPs significantly associated with each trait in our population (*SI Appendix*, Fig. S6). Regardless of trait, no significant SNP was detected using the additive model (Fig. 2*C* and *SI Appendix*, Figs. S9*A*, S10*A*, and S11*A*) after multiple testing correction (Bonferroni threshold, $P < 3 \times 10^{-7}$), consistent with the low power observed in our simulations. The overdominant LMM was fitted to both the trait means and MPH. Significant SNPs were detected for four traits (Fig. 2 *D* and *E* and *SI Appendix*, Figs. S9 *B* and *C*, S10 *B* and *C*, and S11 *B* and *C*), with many more associations for MPH than for the simple trait means (35 vs. 5 significant sites) (*SI Appendix*, Table S3). Significant SNPs collapsed into nine regions; four of these were significant for multiple traits (*SI Appendix*, Table S4). To account for multiple testing across model–phenotypic component combinations, a more stringent Bonferroni correction was applied within each trait [individual significance thresholds (0.05) divided by GWA studies per trait (3) and SNPs (204,753)] (*SI Appendix*, Tables S3 and S4). Most SNPs were significant even after correction within trait (*SI Appendix*, Table S4). Traits with lower heritability, particularly rosette traits of young plants extracted from images, showed no association with any position in the genome. We also did not find any associations with adult rosette size, despite a broad-sense heritability of 0.62. In conclusion, we identified a number of genomic positions that are associated with both the trait means and MPH, suggesting that within-locus interactions, either dominance or overdominance, contribute significantly to nonadditive genetic variance.

**Significant SNPs Contribute Heavily to Genetic Variance.** To assess the contribution of within-locus interactions to genetic variance, the variance explained by each significant SNP was quantified and compared with the variance explained by all tested SNPs. Variance explained by all tested SNPs was computed using a LMM that fitted the appropriate kinship matrix to the trait of



**Fig. 2.** Power analyses of GWA models and a case study. Power of the (*A*) additive and (*B*) overdominant models was calculated for various minor allele frequencies (MAF) and effect sizes. The 0.05, 0.10, and 0.20 curves in *A* all have the value 0.0. The 0.40, 0.60, and 0.80 curves in *B* all have the value 1.0. Variance explained by the focal SNP is considered a proxy for effect size and varies between 5 and 80% of the variance (*Materials and Methods*). Simulations were performed 1,000 times for each combination of MAF and percent of explained variance. Mean power of the simulations as well as the SEM are plotted. $\beta$ is the rate of type II error. Simulations were performed at a type I error rate of 0.05. (*C–E*) Genome-wide *P* values of each test statistic are shown for each association study performed on the DTF trait. Horizontal dashed lines correspond to 5% significance thresholds for both within- ($P < 2.4 \times 10^{-7}$; orange) and across-trait ($P < 8.1 \times 10^{-8}$; blue) Bonferroni corrections. Genomic controls were estimated as the deviation of the observed median test statistics from expected median test statistic. (*C*) Additive model, mean phenotype. (*D*) Overdominant model, mean phenotype. (*E*) Overdominant model, MPH.

interest using a cross-validation strategy. The model was trained with a dataset consisting of 90% of hybrids and then used to predict phenotypes in the test dataset, the remaining 10% hybrids, with 1,000 repetitions. The variance explained by all tested SNPs accounted for 7–56% of the total genetic variance (Fig. 3A) using the additive encoding, whereas it ranged from 18 to 45% (Fig. 3B) using the overdominant encoding. We conclude that our strategy enabled excellent phenotypic predictions and subsequent estimation of variance in our diallel, but we note that variance estimates cannot be necessarily extrapolated to other genotypes (49).

We further estimated the contribution of the significant loci to total phenotypic variation. We found that individual significant SNP generally had a large marginal effect and could explain from 0.02 to 19.6% of the phenotypic variance (Fig. 3C). The contribution of all significant SNPs was calculated using a ridge regression model, together with the cross-validation strategy described above, to account for nonindependence, or linkage, between significant SNPs. Significant SNPs explained up to 20% of the total genetic variance for some traits (LTF MPH and rosette dry mass MPH) (Fig. 3D).

**Multilocus SNP Associations.** In addition to single SNP tests, we searched for multilocus associations using a network-guided approach implemented in the GWA method SConES (50). This approach leverages the protein interaction network of *A. thaliana*



**Fig. 3.** Variance explained by tested SNPs and the characteristics of associated SNPs. (A and B) Variance explained using all tested SNPs was calculated using a cross-validation approach. Mean variance explained and the SEM (1,000 training sets) are plotted for the training (90%) and test (10%) sets. (A) Phenotypic variance of the mean trait value explained by the additive model. (B) Phenotypic variance of MPH explained by the overdominant model. Models that are evaluated using their own training data tend to overfit, hence the values that are close to, or equal to, 1. Variance explained by (C) individual SNPs and (D) all significantly associated SNPs calculated using a cross-validation approach as described for A and B (*Materials and Methods*). (C) Boxplot shows the median (white circle), upper and lower quartiles (black box), and 1.5× the interquartile range (black lines). (D) Mean variance explained and the SEM (1,000 training sets) for the training (90%) and test (10%) sets. (E) MAF distributions for significant SNPs found via single-locus GWA studies. MAFs of all tested SNPs (black) were calculated for an expanded panel of 80 accessions. Random sampling of equal sample sizes (gray) of the test SNPs was repeated 1,000 times. Median values of random draws never reached the median values of the actual data (permutation test, *P* < 0.002).
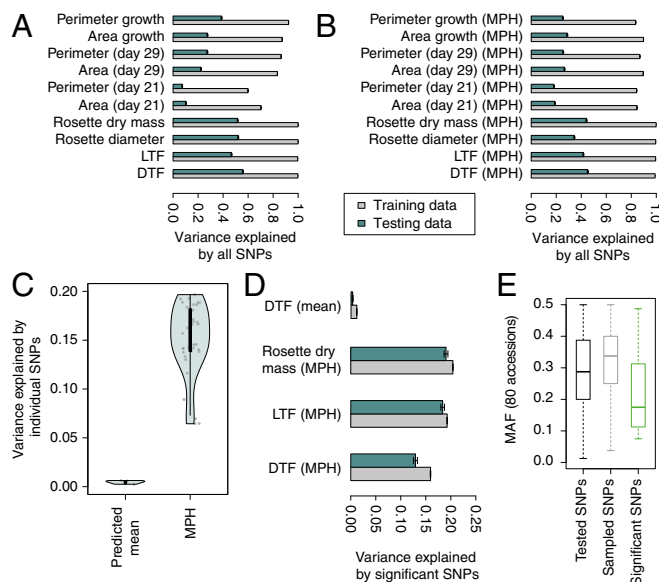
to search for SNPs that together influence a phenotype; however, it does not explicitly test for epistasis between pairs of loci (50). Associated SConES SNPs likely contribute to phenotypic variance either via the sum of multilocus additive effects or allelic heterogeneity at a single locus, where multiple, unlinked SNPs in or near to a gene have similar phenotypic consequences. For each trait investigated with SConES, between 0 and 324 SNPs were linked to the trait of interest (*SI Appendix*, Tables S5 and S6), explaining up to 40% of the total genetic variance of MPH when fitting the overdominant model (*SI Appendix*, Fig. S12A). Less genetic variance could be explained when fitting the overdominant model to the predicted mean phenotype (*SI Appendix*, Fig. S12B). The genetic variance explained by SConES SNPs is not necessarily independent from the variance explained by SNPs detected via traditional GWA mapping, but in this case, only a single SNP was detected with both methods (*SI Appendix*, Tables S4 and S6).

**Significant SNPs and Established Hypotheses for Heterosis.** Established hypotheses regarding the genetic basis of heterosis make specific predictions regarding the allele frequencies of causal loci. Under the dominance hypothesis, causal loci are expected to be rare in the population (reviewed in ref. 8). Minor allele frequencies of all tested SNPs were estimated in the 80 resequenced genomes, from which the 30 diallel parents were drawn (36), and compared with minor allele frequencies of SNPs with significant phenotypic associations either based on single sites (Fig. 3E) or via SConES (*SI Appendix*, Fig. S12C). Significant SNPs had much lower minor allele frequencies than background SNPs (Fig. 3E and *SI Appendix*, Fig. S12C), even though low-frequency variants (<10%) had been removed. Additionally, random sampling of background SNPs demonstrated that their median allele frequencies were always higher than those of the significant SNPs (permutation test, *P* < 0.002 for 1,000 permutations) (Fig. 3E and *SI Appendix*, Fig. S12C). Because of statistical limitations in GWA studies, it is important to note that we were unable to query the effects of truly rare variants.

The significant SNPs were collapsed into nine distinct genomic regions, which included both overdominant and dominant effects (*SI Appendix*, Table S7 and Fig. S13). The most sensitive GWA studies, where MPH was fitted to an overdominant model, detected SNPs in most regions (*SI Appendix*, Tables S3 and S4). Of the nine regions, four behaved dominantly and three behaved overdominantly or pseudooverdominantly with respect to the trait mean (*SI Appendix*, Table S7 and Fig. S13). The remaining two regions tended toward overdominant behavior, but the effect was mild.

If overdominant traits were, in fact, the result of multiple dominant loci, then the magnitude of MPH should increase upon inclusion of additional loci. To test this, multilocus genotypes of dominant regions were correlated with trait means. The phenotypic behavior varied by multilocus genotype and for some combinations did, in fact, exhibit a trend toward overdominance (*SI Appendix*, Fig. S14), suggesting that at least a portion of heterotic phenotypes can be attributed to the combination of multiple, unlinked dominant loci.

The dominance hypothesis predicts that the degree of heterosis in a hybrid will correlate with genetic distance between the parents (51–53). Pairwise genetic distances were calculated across the entire genome at a variety of annotated sites (intergenic, intron, synonymous sites, nonsynonymous sites, etc.). To compare the degree of heterosis across hybrids, MPH values were normalized by the additive component, *a*, or half the distance between the two parental trait values of each cross. Regardless of annotation, correlations between genetic distance and estimates of MPH/*a* were only occasionally significant (*SI Appendix*, Fig. S15), with the direction varying by trait. Flowering time (DTF and LTF) was significantly positively correlated with genetic distance for several polymorphism categories (maximum Spearman rank correlation coefficient <0.16), whereas rosette

perimeter and diameter were negatively correlated with genetic distance (maximum Spearman rank correlation coefficient <|−0.16|). Contrasting directions of heterosis for different traits have been observed before (54), and other studies have also failed to detect a relationship between genetic distance and the magnitude of heterosis (37, 55–57).

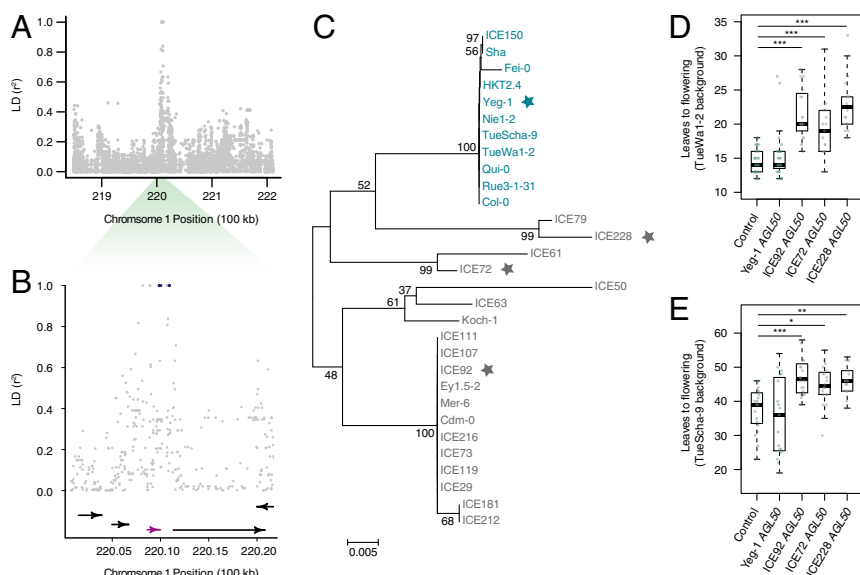**Candidate Genes Are Associated with Relevant Biological Processes.** To gain insight into the biological relevance of each GWA study, we asked whether the top 1,000 SNPs associated with each trait were enriched for specific gene ontology (GO) categories. As expected, flowering time related traits were associated with long-day photoperiodism and photomorphogenesis and, in addition, with GO terms related to posttranscriptional regulation (*SI Appendix*, Table S8). Growth-related traits extracted from the images of young plants were associated with energy production via oxidative phosphorylation in the mitochondria. Although many of these SNPs were not significant using a Bonferroni significance cutoff, the enrichments observed in GO analyses suggest that our study detected additional contributing loci (*SI Appendix*, Table S8).

We measured LD surrounding high-confidence SNPs to identify potential candidate genes underlying the observed heterotic effects. The nine significant regions collapsed into eight linkage blocks (*SI Appendix*, Table S7). In some cases, LD decayed quickly around the significant SNPs, allowing the identification of high-confidence candidate genes (Fig. 4A and *SI Appendix*, Fig. S16 and Table S7). One region in particular, HV1.3, which had a dominant phenotypic effect on leaf number at flowering, exhibits rapid LD decay (Fig. 4 A and B). This short haplotype block spans a single gene, *AGAMOUS-LIKE 50* (*AGL50*), which encodes a MADS-box transcription factor. Many other members of the MADS-box family play critical roles in flowering time control (58–60), but the functions of *AGL50* and its closest paralog *AGL49*, which belong to a poorly characterized clade of the MADS-box family (61), have been unknown.

To examine whether *AGL50* haplotype diversity is predicted by the two associated SNPs located in the downstream intergenic region (Fig. 4B), we sequenced the coding region of *AGL50* from each of the 30 parental accessions (*SI Appendix*, Fig. S17). The phylogeny derived from the inferred amino acid sequences identified two groups distinguished by the diagnostic SNPs (Fig. 4C). Accessions carrying the reference haplotype (blue) flowered on average with 10 fewer leaves than the nonreference haplotype (gray) (mean LTF for nonreference = 35.1 and reference = 24.9; Wilcoxon rank sum test, $P = 0.02$). To directly demonstrate that *AGL50* haplotypes affect flowering time in our population, we transformed genomic fragments of four *AGL50* alleles into two genetic backgrounds, TueWa1-2 and TueScha-9, both belonging to the reference (blue) haplotype. *AGL50* from Yeg-1, a representative of the reference clade that varies at four noncoding positions relative to TueWa1-2 and Tuescha-9, did not change LTF (Fig. 4 D and E), but three nonreference *AGL50* alleles, from ICE92, ICE228, and ICE72, significantly increased LTF (Fig. 4 D and E). Using heterozygous allelic status in a GWA study, we thus uncovered a flowering time role for a gene that had not been linked before to this trait, either by conventional GWA or by mutant analyses.

## Discussion

**The Contribution of Dominance to Genetic Variance.** There is often a discrepancy between the heritability of a trait and inheritance in families on one hand and the genetic variance explained by loci identified in GWA studies on the other hand. This discrepancy is often referred to as "missing heritability," and the potential causes are a subject of ongoing, intense debate in the field of quantitative genetics (reviewed in refs. 62–65). Proposed explanations include the lack of power to detect loci of small effect (reviewed in refs. 62, 66), the importance of rare variants (reviewed in refs. 62, 64), the contribution of multiple different alleles at the same locus (allelic heterogeneity) (50, 67–69), the change in allelic



**Fig. 4.** Candidate region HV1.3 contains *AGL50*, which is sufficient to alter LTF. LD is plotted for (*A*) 400 kb and (*B*) 20 kb surrounding region HV1.3 (*Materials and Methods*). Region HV1.3 (*A*) is also shown in *SI Appendix*, Fig. S16. Significant SNPs associated with LTF MPH are plotted as blue circles (*SI Appendix*, Table S4). Gene models (TAIR10) are shown in *B*, and *AGL50* (AT1G59810) is highlighted in purple. (*C*) Phylogeny built from translated amino acid sequences of *AGL50* from each parental accession using the neighbor-joining method based on the Kimura 2-parameter model implemented in the software MEGA5 (103). Sequences from Bak-2 were excluded for phylogeny construction due to duplication of *AGL50*. Accessions containing the *AGL50* alleles linked to the reference haplotype (TA) of the top two associated SNPs (Chr1:22009862 and Chr1:22009873) are colored blue. Accessions with the alternative haplotype (GG) are shown in gray. *AGL50* alleles from accessions marked with stars were transformed into two genetic backgrounds, TueWa1-2 and TueScha-9. The distribution of LTF in transgenic lines carrying various alleles of *AGL50* in the (*D*) TueWa1-2 and (*E*) TueScha-9 background (*n* = 17 to *n* = 40). Control refers to the nontransgenic background relevant to each plot. Significance was determined using a Wilcoxon rank sum test. *$P < 0.01$; **$P < 0.001$; ***$P < 0.0001$.

effects across environments (65), and the interactions between or within loci (70).

We have leveraged the power of inbred lines and a carefully chosen intercrossing scheme to measure the contribution of nonadditive inheritance to phenotypic variability. The contribution of allelic interactions to genetic variation is most relevant for outcrossing species, but we reasoned that using a collection of replicable, heterozygous lines was a compelling system for exploring the power of nonadditive GWA models. Although nonlinear models are typically ignored in the analysis of continuous traits (62), a considerable portion of genetic variance in our population is attributable to dominance, and the underlying alleles can be mapped when nonadditivity is explicitly considered. Although our approach does not quantify the contribution of nonadditive inheritance to missing heritability in natural populations of *A. thaliana* (which are mostly inbred), our results do argue for the inclusion of nonadditive models in GWA studies because they reveal loci that would go undetected using standard approaches.

We cannot exclude a role for epistasis, but we found that a single heterozygous position can contribute up to 20% of the genetic variance (Fig. 3C) and that the marginal effect of significant SNPs can result from single-locus dominance. An indirect test for multilocus effects using the network-guided SConES approach (50) suggests that allelic heterogeneity or multilocus additive effects may account for additional phenotypic variation in our population because SConES mostly found other significant SNPs than the single-locus scans. Because neither set of SNPs can explain all of the genetic variance, we hypothesize that both intralocus and interlocus interactions contribute to nonadditivity in our population and that the genetic basis for both contributions is largely nonoverlapping.

**Arabidopsis thaliana in the Context of Traditional Heterosis Hypotheses.** For nearly 100 years, geneticists have sought to develop a unified model explaining heterosis. Three leading hypotheses of intralocus interactions have been developed. The dominance hypothesis suggests that individuals in a population carry a suite of rare, slightly deleterious mutations that have not yet been purged by purifying selection (16–18, 71, 72). The ability of selection to remove weakly deleterious mutations is reduced when effective population sizes are small, which is typical for selfing species. Correspondingly, the dominance hypothesis has generally received the most empirical support from studies of inbreeding depression (73–76). If heterosis is the reverse of inbreeding depression, then the degree of heterosis should positively correlate with the genetic distance between parents, and causal alleles should be rare with small phenotypic effects (8, 51–53, 73). Several dominantly acting loci have been shown to contribute to heterosis (29, 77, 78), and more recently, heterosis-associated loci in maize have been shown to be enriched for deleterious mutations (43). Although dominance remains the prevailing explanation, some assumptions of this hypothesis are not consistently supported; the correlation between the degree of heterosis and the genetic distance between parents is not always evident (37, 55–57, 79–82), and there are loci with moderate effect sizes (29, 77, 78).

The overdominance hypothesis suggests that a very small number of overdominant loci with large effects explains the majority of heterotic phenotypes (3, 4, 9). This alternative hypothesis is good news for breeding programs because a few major-effect loci are much more easily introduced into different backgrounds than a large number of small-effect loci. A number of studies have identified overdominant QTL associated with hybrid vigor (23, 24, 78, 81, 83, 84), but molecular identification of casual variants is rare. Although a few cases of truly overdominant loci have been confirmed (10–13), in some cases, fine mapping of overdominant QTL has separated a single overdominant locus into multiple,

dominant loci acting in repulsion (85, 86), a situation called pseudooverdominance, which represents the third common hypothesis for heterosis.

As discussed above, characteristics of dominant and overdominant alleles underlying heterosis, including their expected effect size and allele frequency, have been predicted by assuming that the trait in question is related to fitness. We have focused on two groups of traits, flowering time and plant growth/size. We did not directly estimate a fitness proxy such as seed set, but we reasoned that traits shown to have adaptive value in *A. thaliana*, such as flowering time (87–89), will be subject to evolutionary forces comparable to those acting directly on fitness. Although flowering time is locally adaptive in *A. thaliana*, the correlation between flowering time and fitness varies by accession and environment (90–92). Large-effect mutations in flowering time genes can significantly perturb fitness, but their effect is not directional and varies by genetic background (90, 93), providing one possible explanation for variable correlation between these two traits.

Our data are a poor fit for the dominance hypothesis because we found overdominant and dominant loci of medium to large effect. The variants at these loci, although segregating at lower frequency than background SNPs, do not classify as rare by population genetic standards. Furthermore, there is only a weak positive correlation of nonadditivity with genetic distance for most traits, and the strongest evidence for any relationship is a negative correlation with rosette diameter (*SI Appendix*, Fig. S15), with the caveat that any correlation based on a large number of small-effect loci may be obscured by the moderate-effect size loci that we detected in the GWA studies. Several previous studies of heterosis using controlled crosses in *A. thaliana* have identified loci that exhibit all possible modes of gene action, including additive, dominant, and epistatic interactions (26–28, 94–96). Lack of support for the major heterosis hypotheses comes also from studies in crop species, particularly in maize and rice. Maize is a classical model for investigating heterosis, and this outcrossing species is cited frequently as supporting the dominance model of heterosis (18), but several overdominant QTL have been identified (29, 30, 81). Rice, in contrast, has been proposed as a system that supports the overdominance hypothesis, but all three modes of gene action have been uncovered in this predominantly selfing species as well (23, 24, 29, 77, 78). The dichotomy between these two model crop species has been attributed to their alternative mating strategies; the large-effect loci that we have found in *A. thaliana*, a selfing species, support this argument.

That different single-locus allelic interactions can underlie heterotic phenotypes in multiple species suggests that both single-gene dominance and overdominance truly occur or that pseudooverdominance is more prevalent than expected. Additionally, it is also possible that putatively deleterious, dominant loci have pleotropic effects. If such loci contribute positively to a second phenotype, they could be retained during evolution for longer than expected from their deleterious effects. Regardless of genetic behavior, the existence of large-effect, Mendelian loci driving heterosis is a considerable boon to plant breeding programs where they could easily be integrated into elite material.

## Materials and Methods

**Generation of Plant Material.** A half diallel was constructed by manually intercrossing 30 inbred strains of *A. thaliana* (36), facilitated by male sterility induced by an artificial miRNA targeting the homeotic genes *AP3* (22). In addition to the 435 hybrid combinations generated with this method, 30 manual self-crosses of the parental strains were performed using the same strategy. This ensured that the maternal environments of the hybrid and parental genotypes were equivalent. A list of hybrid and parental genotypes used in this study can be found in *SI Appendix*, Table S1.

**Experimental Design.** In total, five replicates of 495 genotypes were surveyed in this experiment (435 hybrid genotypes, 30 parental genotypes from

manual crosses, and 30 self-fertilized parental genotypes). Five unsterilized seeds for each replicate were aliquoted into 1.5-mL tubes with 500 μL of ddH₂0. Seeds were stratified in the dark at 4 °C for 10 days. After stratification, seeds were sown into soil (CL T Topferde; www.einheitserde.de) pots in a completely randomized design. Flats were covered with humidity domes and placed into 16 °C growth chambers under long-day conditions (16 hours light: 8 hours dark) at a relative humidity of 65%. Light bulbs were a mixture of Sylvania Cool White Deluxe to Warm White Deluxe fluorescent bulbs (4:2) (www.havells-sylvania.com/en-GB/sylvania). Humidity domes were removed after 1 week, and pots were manually thinned to one plant per pot. Plants were subsequently phenotyped for a variety of traits: days to first open flower (DTF), rosette leaf count at the first open flower (LTF), rosette diameter, and rosette dry mass. Once the plants had produced about 10 siliques, they were harvested, and rosette diameters were measured. The rosettes were placed into paper bags, dried at 80 °C for 24 hours, and weighed. Additionally, images of each tray were taken at days 21 and 29. From these images the following measurements were extracted using a custom ImageJ (97) macro: area (day 21 and 29), perimeter (day 21 and 29), area growth [(day 29 − day 21)/8 d], and perimeter growth [(day 29 − day 21)/8 d]. In brief, the macro automatically segmented the images by removing the background and returned rosette area and perimeter values in pixels for each plant. Because the maternal plants were hemizygous for the artificial miRNA targeting *AP3*, progeny derived from these crosses were segregating for the transgene. Plants with the transgene were easily identified based on their floral and fruit morphology. To ensure that the transgene did not alter the measured phenotypes, we recorded the artificial miRNA status of each plant for use as a covariate in later analyses. Additionally, the dates that the plants were harvested for rosette measurements were recorded for use as a potential covariate.

**Handling of Missing Data and Data Normalization.** Overall, germination rates were high in this experiment. Out of 495 surveyed genotypes, only 9 completely failed to germinate (7 hybrids and 2 manually selfed parents), and these lines were excluded from further analyses (*SI Appendix*, Table S1). Of the remaining lines, 98% of plants germinated. Most germination failures only occurred in a single replicate (*SI Appendix*, Fig. S18). In these cases (58 in total), the missing phenotypes were imputed as the mean of the phenotyped replicates for each genotype. After exclusion of genotypes with failed germination and imputation of the remaining missing data, each phenotype was Box–Cox transformed to improve the normality of the data (*SI Appendix*, Fig. S19).

**Estimation of GCA, SCA, and Heritability.** A traditional ANOVA approach was not appropriate for our data, because there are some missing data (35). Instead, variance components were estimated using a linear mixed model implemented in SAS using a restricted maximum likelihood estimation method (98). The SAS code is available in *SI Appendix*, Text S1, and is a modified version of the code available on Fikret Isik's webpage (www4.ncsu.edu/~fisik/Analysis%20of%20Diallel%20Progeny%20Test%20with%20SAS.pdf). Only hybrid genotypes were used. The following linear mixed model was fitted to the transformed data:

$$Y_{jkl} = \mu + G_j + G_k + S_{jk} + E_{jkl}.$$

Here $Y_{jkl}$ is the *l*th phenotypic observation for the *jk*th cross, $\mu$ is the overall mean, $G_j$ or $G_k$ is the random GCA of the *j*th female or the *k*th male, $S_{jk}$ is the random SCA of the *j*th female and the *k*th male, and $E_{jkl}$ is the error term. All terms were expected to be normally distributed. This model can also be written in matrix format:

$$y = X\beta + Z\gamma + \varepsilon.$$

Here $y$ is a vector of observations, $\beta$ is a vector of the fixed effects parameter (overall mean), $\gamma$ is the vector of random effects parameters (GCA and SCA), $\varepsilon$ is the random error vector, $X$ is the known design matrix for the fixed effects, and $Z$ is the known design matrix for the random effects. In SAS, the $Z$ design matrix was constructed by hand using PROC IML to associate each individual with its respective parents. Next, PROC MIXED was run on the data using the above model. Variance components and covariances of variance components were extracted from the model and used to calculate both broad- and narrow-sense heritability (as well as their SEs). Because our parents were not derived from a randomly mated population, the additive ($\sigma_A^2$) and dominance ($\sigma_D^2$) genetic variance and the total phenotypic variance ($\sigma_P^2$) in our data were as follows (35):

$$\sigma_A^2 = 2\sigma_{GCA}^2,$$

$$\sigma_D^2 = \sigma_{SCA}^2,$$

$$\sigma_P^2 = 2\sigma_{GCA}^2 + \sigma_{SCA}^2 + \sigma_{Error}^2.$$

Both narrow- ($h_n^2$) and broad-sense ($H_b^2$) heritabilities were calculated from these values (35):

$$H_b^2 = \frac{\sigma_A^2 + \sigma_D^2}{\sigma_P^2},$$

$$h_n^2 = \frac{\sigma_A^2}{\sigma_P^2}.$$

**Estimation of Mean Genotypic Values.** A linear mixed model was fitted to each Box–Cox transformed phenotype using the package lme4 in the R statistical framework (99). For each phenotype, the following model was fitted:

$$Y_{jkl} = G_{jk} + A_{jk} + E_{jkl},$$

where $G_{jk}$ is the random genotypic effect of the *j*th female and the *k*th male, $A_{jk}$ is the random effect of the *amiR AP3* transgene on the hybrid cross of *j*th female and the *k*th male, and $E_{jkl}$ is the error term. For each phenotype, the above model was fitted with and without the transgene variable, and the significance of this term was tested. In a few cases the transgene term was not significant and was subsequently removed from the model (DTF, LTF, and dry mass). In the remaining cases, the transgene explained only 0.02–2.18% of the total phenotypic variance. After model fitting, the coefficients of each genotype were extracted from the model and used for all subsequent analyses. Broad-sense heritability was also calculated from these models:

$$H_b^2 = \frac{\sigma_G^2}{\sigma_P^2}.$$

Here $\sigma_G^2$ is the variance due to the genotype term ($G_{jk}$), and $\sigma_P^2$ is the total phenotypic variance. The broad-sense heritability estimates were comparable to those derived from SAS (Fig. 1).

**Calculation of Midparent Heterosis.** The predicted genotypic values of hybrid and manually selfed parental genotypes were extracted from the linear model. Using these values, standard quantitative genetic components of phenotype were calculated (7). MPH was calculated as the distance of the hybrid phenotype from the midparent value, or mean of the two parental genotypes. Two of the 30 manually selfed parental genotypes did not germinate, and as a result, MPH could not be calculated for hybrids generated from these two parents (Bak-2 and ICE61) (*SI Appendix*, Table S1). For subsequent analyses, only 372 hybrid genotypes were used.

**Generation of F₁ Genotypes and SNP Filtering.** In silico F₁ genotypes were constructed from parental genome sequences (36) that had been filtered to remove (*i*) all sites that lacked complete information, (*ii*) all sites that were not polymorphic, (*iii*) all triallelic sites (with respect to the reference), and (*iv*) all singletons. After filtering, 723,403 SNPs remained. Because the parental genotypes are few, there is extensive long-distance LD between sites. To remove such sites, we first encoded all 723,402 SNPs using the standard additive 0,1,2 encoding, where 0 is the major, 1 is the heterozygous, and 2 is the minor allele. After encoding, 75,346 SNPs were only observed once within this population. We then created categories for how often a specific SNP pattern across all individuals was observed within our dataset (pattern occurrence). These categories ranged from 2 to 7,364. For example, a SNP is located in category 2 if this SNP shares the same pattern with exactly one other SNP in the genome and in category 1,000 if the SNP in question shares the same pattern with exactly 999 other SNPs. In *SI Appendix*, Fig. S7, the cumulative number of SNPs for all categories is plotted. We observe that 32.97% of all our SNPs fall into the categories 1–10, which includes all distinct SNPs plus the number of SNPs for each of the categories from 2 to 10. Approximately 38% of all SNPs fall into the categories 100–7,364. Next we evaluated whether SNPs with shared patterns were located on the same chromosome or distributed across multiple chromosomes. *SI Appendix*, Fig. S8, shows the distribution of SNPs across chromosomes for categories 2–20. For the final SNP set, we allowed SNPs to share their pattern with up to nine other positions (categories 1–10), but we removed all sites that exhibited complete long-distance LD across chromosomes. The final dataset consisted of 204,753 SNPs, and these sites were used for all association mapping studies.

**Genome-Wide Association Mapping (Additive Model).** All GWA analyses were conducted using the easyGWAS framework (46). We used a local copy of easyGWAS and custom C/C++ and Python implementations of the FaSTLMM (47) algorithm. For the additive model, the homozygous major allele is encoded with 0, the heterozygous genotype with 1, and the homozygous minor allele with 2. The genetic similarity between all genotypes was estimated by computing the realized relationship kinship matrix (48) on the additively encoded genotype data. This kinship matrix was used in the FaSTLMM model to account for confounding due to population stratification and cryptic relatedness. The additive model was only run on the predicted phenotypic values. Genomic control (GC) values were computed to assess the degree of inflated test statistics (100). GC is measuring the deviation of the observed median test statistics from the expected one. GC values larger than 1 are indicative of inflated P values, whereas values smaller than 1 are indicative of deflated P values. GC values for each GWA can be found in *SI Appendix*, Table S9, and QQ plots for traits with significant SNPs can be found in *SI Appendix*, Fig. S20.

**Genome-Wide Association Mapping (Overdominant Model).** We conducted GWA analyses with an overdominant genotype encoding, where both the homozygous minor and homozygous major alleles are encoded as 0 and the heterozygous genotype is encoded as 1. The kinship matrix was computed on the overdominantly encoded data. Using the overdominant encoding, GWA mapping was performed on both the predicted phenotypic values of the hybrids as well as MPH of each strain.

**Multiple Hypothesis Testing Correction.** To account for multiple hypothesis testing, we used a conservative 5% Bonferroni threshold of 0.05/[number of tested SNPs (204,753)] = $2.4 \times 10^{-7}$. This correction was performed within each study, and significant results are reported in *SI Appendix*, Tables S3 and S4. Additionally, we performed an even more stringent correction by accounting for the number of GWA analyses per trait (3). In this case the Bonferroni threshold was equal to $2.4 \times 10^{-7}/3 = 8.1 \times 10^{-8}$. The results from this test correction are reported in *SI Appendix*, Tables S3 and S4.

**Estimation of Variance Explained by All SNPs.** We computed how much of the phenotypic variance could be attributed to the genetic contribution (random effect) using a cross-validation approach. We generated 1,000 randomly drawn training sets (containing 90% of all hybrid genotypes) and testing sets (remaining 10% of genotypes). We then trained the LMM using only the kinship matrix (random effect) on the training data and subsequently predicted the phenotype $\hat{y}$ of the remaining testing set. Predictions were obtained as follows:

$$\hat{y} = C_{\text{test}}\tilde{\beta} + K_{\text{test}}\left(K_{\text{train}} + \tilde{\delta}I\right)^{-1}\left(y_{\text{train}} - C_{\text{train}}\tilde{\beta}\right),$$

where $C$ is a vector of ones (or different covariates if given), $K$ is the kinship matrix, and $\tilde{\beta}$ and $\tilde{\delta}$ are the estimated parameters from the training step of the LMM. We then computed variance explained as follows:

$$v(y_{\text{test}}, \hat{y}) = 1 - \frac{\text{var}(y_{\text{test}} - \hat{y})}{\text{var}(y_{\text{test}})},$$

where var() is the variance. Note that this measure might be negative, and in such cases the phenotypic mean would provide a better fit than the actual trained model. Results were averaged across all 1,000 training sets.

**Estimation of Variance Explained by Individual Significant SNPs.** Next, we estimated the variance explained by each individual SNP. For each significantly associated SNP we generated 1,000 randomly drawn training sets (containing 90% of all hybrid genotypes) and testing sets (remaining 10% of lines). We then computed variance explained by a single SNP by fitting a linear regression:

$$v_{\text{SNP}} = 1 - \frac{\text{var}\left(y_{\text{test}} - \tilde{\beta}X\right)}{\text{var}(y_{\text{test}})},$$

where $\tilde{\beta}$ is the estimated parameter from the linear model and $X$ is the associated SNP. The parameter $\tilde{\beta}$ is estimated as follows:

$$\tilde{\beta} = \left(X^TX\right)^{-1}X^Ty,$$

where $X^T$ is the transposed matrix $X$. Results were averaged across all 1,000 training sets. Note that these phenotypic predictions are only relevant to the current dataset and that the results (i.e., variance explained by each site) need not generalize to genotypes outside of this dataset.

**Estimation of Variance Explained by All Significant SNPs.** It is important to note that one cannot sum up the variance explained by all individual SNPs to obtain the variance explained by all significantly associated SNPs, because the positions are not entirely independent of one another, predominantly due to LD. To estimate the variance explained by all significantly associated SNPs, we trained a ridge regression on $X$, where $X$ contains all significantly associated SNPs. Ridge regression includes a penalty term to regularize the weight of each SNP and thus implicitly takes the relatedness between individual SNPs into account:

$$\tilde{\beta} = \left(X^TX + \lambda I\right)^{-1}X^Ty,$$

where $\lambda$ is the penalty term. $\lambda$ is optimized by performing an internal line search for a range of $\lambda$ values: $\lambda = \{1e^{-3}, 1e^{-2}, 1e^{-1}, 1, 1e^1, 1e^2, 1e^3\}$. Again, 1,000 cross-validation sets were run for each model.

**Power Analyses.** To evaluate the power of the different encoding strategies, we performed a simulation analysis in which we measured the power of each test with respect to the variance explained by the causal SNP, the minor allele frequency of the causal SNP, and the SNP encoding. The simulations were performed with both the additive and overdominant SNP encodings. We binned the tested SNPs (204,753) according to their minor allele frequency {0.10–0.15,...,0.45–0.50}. As the background covariance matrix (kinship matrix) we used the realized relationship matrix based on all SNPs (204,753), applying the appropriate encoding (48). For combinations of factors (variance explained, minor allele frequency, and SNP encoding), we first randomly chose a causal SNP with the selected minor allele frequency from our genotypic data. We simulate the phenotype as

$$y = X\beta + \epsilon,$$

where $X$ is the causal SNP and $\beta$ is the regression coefficient. We let the proportion of variance explained by the SNP vary between (0.05, 0.10, 0.20, 0.40, 0.60, 0.80). The remaining variance is explained by Gaussian distributed noise:

$$\epsilon \sim N(0, 1 - v_{\text{SNP}}I),$$

where $v_{\text{SNP}}$ is the variance explained by the focal SNP. Each combination of factors (variance explained, minor allele frequency, and SNP encoding) was repeated 1,000 times. Results show the power, or 1 minus the probability of not detecting the causal SNP, averaged over all repetitions, along with the SEs (Fig. 2 *A* and *B*). These simulations were performed using a type I error rate of 0.05.

**Characterization of Significant Peaks and Identification of Candidate Genes.** For GWA studies, we considered only SNPs with complete genotype information in our parental panel, but this approach removes some potentially relevant polymorphism. To characterize the decay of linkage disequilibrium (LD) around peaks and to develop a candidate gene list, we used a less stringent cutoff of 70% complete information in all sites and identified additional candidate SNPs based on linkage to significant sites using the program PLINK 1.9 (101). SNPs within 200 kb of a significant SNP, in complete LD ($r^2 = 1$), and with a minor allele frequency greater than 0.1 were collapsed into the eight candidate regions listed in *SI Appendix*, Table S7. Two peaks on chromosome 3 were collapsed into a single large region using this approach because of their physical proximity and extended LD in this region. The regional information was used to develop candidate lists. Decay of LD around these peaks was calculated from the reference SNP in each of the above-described regional LD blocks for up to 200 kb surrounding the focal SNP (Fig. 4 *A* and *B* and *SI Appendix*, Fig. S16).

**Gene Ontology Analysis.** For each GWA analysis using the overdominant SNP encoding, the top 1,000 most associated SNPs were compared against the complete set of tested SNPs using the SNP2GO library in the R statistical computing environment (102). Significance was established using a 5% false discovery rate threshold within each trait.

**Transgenic Analysis of *AGL50*.** A 1.8-kb genomic fragment (Chr1:22,007,995-22,009,824) was PCR amplified using primers TTGGGAAGATCGATTGTCTCTTA and CTTGCAAGCTCACTGATAGAAAGT from genomic DNA of each parental accession and Sanger sequenced. The *AGL50* fragments from Yeg-1, ICE228, ICE72, and ICE92 were cloned into the binary vector pMLBart for plant transformation. $T_1$ transformants were selected on soil treated with 0.1% glufosinate (Basta) and transferred to individual pots for phenotyping at 2 weeks after sowing. *AGL50* sequences from 29 accessions are available at GenBank under accession nos. KX581758–KX581786.

1. Darwin C (1876) *The Effects of Cross and Self Fertilisation in the Vegetable Kingdom* (John Murray, London), pp viii, 482.
2. Shull GH (1914) Duplicate genes for capsule-form in *Bursa bursa-pastoris*. *Z Vererbungsl* 12(1):97–149.
3. Shull GH (1908) The composition of a field of maize. *J Hered* os-4(1):296–301.
4. East EM (1908) Inbreeding in corn. *Rep Conn Agric Exp Stn* 1907:419–428.
5. Duvick DN (2001) Biotechnology in the 1930s: The development of hybrid maize. *Nat Rev Genet* 2(1):69–74.
6. Riedelsheimer C, et al. (2012) Genomic and metabolic prediction of complex heterotic traits in hybrid maize. *Nat Genet* 44(2):217–220.
7. Falconer DS, Mackay TFC (1996) *Introduction to Quantitative Genetics* (Addison Wesley Longman, Harlow, UK), 4th Ed.
8. Charlesworth D, Willis JH (2009) The genetics of inbreeding depression. *Nat Rev Genet* 10(11):783–796.
9. Schwartz D, Laughner WJ (1969) A molecular basis for heterosis. *Science* 166(3905):626–627.
10. Rédei GP (1962) Single locus heterosis. *Z Vererbungsl* 93(1):164–170.
11. Vrebalov J, et al. (2002) A MADS-box gene necessary for fruit ripening at the tomato *ripening-inhibitor (rin)* locus. *Science* 296(5566):343–346.
12. Krieger U, Lippman ZB, Zamir D (2010) The flowering gene *SINGLE FLOWER TRUSS* drives heterosis for yield in tomato. *Nat Genet* 42(5):459–463.
13. Guo M, et al. (2014) Maize *ARGOS1 (ZAR1)* transgenic alleles increase hybrid maize yield. *J Exp Bot* 65(1):249–260.
14. Todesco M, et al. (2014) Activation of the *Arabidopsis thaliana* immune system by combinations of common *ACD6* alleles. *PLoS Genet* 10(7):e1004459.
15. Smith LM, Bomblies K, Weigel D (2011) Complex evolutionary events at a tandem cluster of *Arabidopsis thaliana* genes resulting in a single-locus genetic incompatibility. *PLoS Genet* 7(7):e1002164.
16. Davenport CB (1908) Degeneration, albinism and inbreeding. *Science* 28(718):454–455.
17. Bruce AB (1910) The Mendelian theory of heredity and the augmentation of vigor. *Science* 32(827):627–628.
18. Jones DF (1917) Dominance of linked factors as a means of accounting for heterosis. *Proc Natl Acad Sci USA* 3(4):310–312.
19. Bomblies K, et al. (2007) Autoimmune response as a mechanism for a Dobzhansky-Muller-type incompatibility syndrome in plants. *PLoS Biol* 5(9):e236.
20. Alcázar R, García AV, Parker JE, Reymond M (2009) Incremental steps toward incompatibility revealed by *Arabidopsis* epistatic interactions modulating salicylic acid pathway activation. *Proc Natl Acad Sci USA* 106(1):334–339.
21. Alcázar R, et al. (2010) Natural variation at Strubbelig Receptor Kinase 3 drives immune-triggered incompatibilities between *Arabidopsis thaliana* accessions. *Nat Genet* 42(12):1135–1139.
22. Chae E, et al. (2014) Species-wide genetic incompatibility analysis identifies immune genes as hot spots of deleterious epistasis. *Cell* 159(6):1341–1351.
23. Li Z, Pinson SRM, Park WD, Paterson AH, Stansel JW (1997) Epistasis for three grain yield components in rice (*Oryza sativa* L.). *Genetics* 145(2):453–465.
24. Li ZK, et al. (2001) Overdominant epistatic loci are the primary genetic basis of inbreeding depression and heterosis in rice. I. Biomass and grain yield. *Genetics* 158(4):1737–1753.
25. Luo X, et al. (2009) Additive and over-dominant effects resulting from epistatic loci are the primary genetic basis of heterosis in rice. *J Integr Plant Biol* 51(4):393–408.
26. Kusterer B, et al. (2007) Heterosis for biomass-related traits in *Arabidopsis* investigated by quantitative trait loci analysis of the triple testcross design with recombinant inbred lines. *Genetics* 177(3):1839–1850.
27. Kusterer B, et al. (2007) Analysis of a triple testcross design with recombinant inbred lines reveals a significant role of epistasis in heterosis for biomass-related traits in *Arabidopsis*. *Genetics* 175(4):2009–2017.
28. Melchinger AE, et al. (2007) Genetic basis of heterosis for growth-related traits in *Arabidopsis* investigated by testcross progenies of near-isogenic lines reveals a significant role of epistasis. *Genetics* 177(3):1827–1837.
29. Garcia AA, Wang S, Melchinger AE, Zeng ZB (2008) Quantitative trait loci mapping and the genetic basis of heterosis in maize and rice. *Genetics* 180(3):1707–1724.
30. Guo T, et al. (2014) Genetic basis of grain yield heterosis in an "immortalized F₂" maize population. *Theor Appl Genet* 127(10):2149–2158.
31. Atwell S, et al. (2010) Genome-wide association study of 107 phenotypes in *Arabidopsis thaliana* inbred lines. *Nature* 465(7298):627–631.
32. Horton MW, et al. (2012) Genome-wide patterns of genetic variation in worldwide Arabidopsis thaliana accessions from the RegMap panel. *Nat Genet* 44(2):212–216.
33. Consortium G; 1001 Genomes Consortium. Electronic address: magnus.nordborg@gmi.oeaw.ac.at; 1001 Genomes Consortium (2016) 1,135 genomes reveal the global pattern of polymorphism in *Arabidopsis thaliana*. *Cell* 166(2):481–491.
34. Bomblies K, et al. (2010) Local-scale patterns of genetic variability, outcrossing, and spatial structure in natural stands of *Arabidopsis thaliana*. *PLoS Genet* 6(3):e1000890.
35. Lynch M, Walsh B (1998) *Genetics and Analysis of Quantitative Traits* (Sinauer, Sunderland, MA), pp xvi, 980.
36. Cao J, et al. (2011) Whole-genome sequencing of multiple *Arabidopsis thaliana* populations. *Nat Genet* 43(10):956–963.
37. Meyer RC, Törjék O, Becher M, Altmann T (2004) Heterosis of biomass production in Arabidopsis. Establishment during early development. *Plant Physiol* 134(4):1813–1823.
38. Comstock RE, Robinson HF (1948) The components of genetic variance in populations of biparental progenies and their use in estimating the average degree of dominance. *Biometrics* 4(4):254–266.
39. Sasieni PD (1997) From genotypes to genes: Doubling the sample size. *Biometrics* 53(4):1253–1261.
40. Sladek R, et al. (2007) A genome-wide association study identifies novel risk loci for type 2 diabetes. *Nature* 445(7130):881–885.
41. Zhang Q, et al. (1994) A diallel analysis of heterosis in elite hybrid rice based on RFLPs and microsatellites. *Theor Appl Genet* 89(2-3):185–192.
42. Ben-Israel I, Kilian B, Nida H, Fridman E (2012) Heterotic trait locus (HTL) mapping identifies intra-locus interactions that underlie reproductive hybrid vigor in Sorghum bicolor. *PLoS One* 7(6):e38993.
43. Mezmouk S, Ross-Ibarra J (2014) The pattern and distribution of deleterious mutations in maize. *G3 (Bethesda)* 4(1):163–171.
44. Laiba E, Glikaite I, Levy Y, Pasternak Z, Fridman E (2016) Genome scan for nonadditive heterotic trait loci reveals mainly underdominant effects in *Saccharomyces cerevisiae*. *Genome* 59(4):231–242.
45. Shapira R, David L (2016) Genes with a combination of over-dominant and epistatic effects underlie heterosis in growth of *Saccharomyces cerevisiae* at high temperature. *Front Genet* 7:72.
46. Grimm D, et al. (2012) easyGWAS: An integrated interspecies platform for performing genome-wide association studies. arXiv:1212.4788v1211.
47. Lippert C, et al. (2011) FaST linear mixed models for genome-wide association studies. *Nat Methods* 8(10):833–835.
48. Hayes BJ, Visscher PM, Goddard ME (2009) Increased accuracy of artificial selection by using the realized relationship matrix. *Genet Res* 91(1):47–60.
49. Wray NR, et al. (2013) Pitfalls of predicting complex traits from SNPs. *Nat Rev Genet* 14(7):507–515.
50. Azencott CA, Grimm D, Sugiyama M, Kawahara Y, Borgwardt KM (2013) Efficient network-guided multi-locus association mapping with graph cuts. *Bioinformatics* 29(13):i171–i179.
51. Charcosset A, Lefort-Buson M, Gallais A (1991) Relationship between heterosis and heterozygosity at marker loci: A theoretical computation. *Theor Appl Genet* 81(5):571–575.
52. Charcosset A, Essioux L (1994) The effect of population structure on the relationship between heterosis and heterozygosity at marker loci. *Theor Appl Genet* 89(2-3):336–343.
53. Bernardo R (1992) Relationship between single-cross performance and molecular marker heterozygosity. *Theor Appl Genet* 83(5):628–634.
54. Flint-Garcia SA, Buckler ES, Tiffin P, Ersoz E, Springer NM (2009) Heterosis is prevalent for multiple traits in diverse maize germplasm. *PLoS One* 4(10):e7433.
55. Cerna FJ, Cianzio SR, Rafalski A, Tingey S, Dyer D (1997) Relationship between seed yield heterosis and molecular marker heterozygosity in soybean. *Theor Appl Genet* 95(3):460–467.
56. Liu ZQ, Pei Y, Pu ZJ (1999) Relationship between hybrid performance and genetic diversity based on RAPD markers in wheat, *Triticum aestivum* L. *Plant Breed* 118(2):119–123.
57. Riday H, Brummer EC, Campbell TA, Luth D, Cazcarro PM (2003) Comparisons of genetic and morphological distance with heterosis between *Medicago sativa* subsp. *sativa* and subsp. *falcata*. *Euphytica* 131(1):37–45.
58. Smaczniak C, Immink RG, Angenent GC, Kaufmann K (2012) Developmental and evolutionary diversity of plant MADS-domain factors: Insights from recent studies. *Development* 139(17):3081–3098.
59. Posé D, et al. (2013) Temperature-dependent regulation of flowering by antagonistic FLM variants. *Nature* 503(7476):414–417.
60. Lee JH, et al. (2013) Regulation of temperature-responsive flowering by MADS-box transcription factor repressors. *Science* 342(6158):628–632.
61. Parenicová L, et al. (2003) Molecular and phylogenetic analyses of the complete MADS-box transcription factor family in *Arabidopsis*: New openings to the MADS world. *Plant Cell* 15(7):1538–1551.
62. Manolio TA, et al. (2009) Finding the missing heritability of complex diseases. *Nature* 461(7265):747–753.
63. Eichler EE, et al. (2010) Missing heritability and strategies for finding the underlying causes of complex disease. *Nat Rev Genet* 11(6):446–450.
64. Gibson G (2012) Rare and common variants: Twenty arguments. *Nat Rev Genet* 13(2):135–145.
65. Mackay TF (2014) Epistasis and quantitative traits: Using model organisms to study gene-gene interactions. *Nat Rev Genet* 15(1):22–33.
66. Ehrenreich IM, et al. (2010) Dissection of genetically complex traits with extremely large pools of yeast segregants. *Nature* 464(7291):1039–1042.
67. Sugiyama M, Azencott C-A, Grimm D, Kawahara Y, Borgwardt KM (2014) Multi-task feature selection on multiple networks via maximum flows. *Proc SIAM Int Conf Data Min* 2014:199–207.
68. Llinares-López F, et al. (2015) Genome-wide detection of intervals of genetic heterogeneity associated with complex traits. *Bioinformatics* 31(12):i240–i249.
69. Sham PC, Purcell SM (2014) Statistical power and significance testing in large-scale genetic studies. *Nat Rev Genet* 15(5):335–346.
70. Bloom JS, Ehrenreich IM, Loo WT, Lite TL, Kruglyak L (2013) Finding the sources of missing heritability in a yeast cross. *Nature* 494(7436):234–237.
71. Fisher RA (1930) *The Genetical Theory of Natural Selection* (Oxford University Press, Oxford), p 318.
72. Kimura M (1983) Rare variant alleles in the light of the neutral theory. *Mol Biol Evol* 1(1):84–93.

73. Charlesworth D, Charlesworth B (1987) Inbreeding depression and its evolutionary consequences. *Annu Rev Ecol Syst* 18:237–268.
74. Barrett SCH, Charlesworth D (1991) Effects of a change in the level of inbreeding on the genetic load. *Nature* 352(6335):522–524.
75. Willis JH (1992) Genetic analysis of inbreeding depression caused by chlorophyll-deficient lethals in *Mimulus guttatus*. *Heredity* 69(6):562–572.
76. Crow JF (1993) Mutation, mean fitness, and genetic load. *Oxford Surv Evol Biol* 9:3–42.
77. Xiao J, Li J, Yuan L, Tanksley SD (1995) Dominance is the major genetic basis of heterosis in rice as revealed by QTL analysis using molecular markers. *Genetics* 140(2):745–754.
78. Hua J, et al. (2003) Single-locus heterotic effects and dominance by dominance interactions can adequately explain the genetic basis of heterosis in an elite rice hybrid. *Proc Natl Acad Sci USA* 100(5):2574–2579.
79. Smith OS, Smith JSC, Bowen SL, Tenborg RA, Wall SJ (1990) Similarities among a group of elite maize inbreds as measured by pedigree, F$_1$ grain yield, grain yield, heterosis, and RFLPs. *Theor Appl Genet* 80(6):833–840.
80. Barbosa AMM, et al. (2003) Relationship of intra- and interpopulation tropical maize single cross hybrid performance and genetic distances computed from AFLP and SSR markers. *Euphytica* 130(1):87–99.
81. Stuber CW, Lincoln SE, Wolff DW, Helentjaris T, Lander ES (1992) Identification of genetic factors contributing to heterosis in a hybrid from two elite maize inbred lines using molecular markers. *Genetics* 132(3):823–839.
82. Larièpe A, et al. (2012) The genetic basis of heterosis: Multiparental quantitative trait loci mapping reveals contrasted levels of apparent overdominance among traits of agronomical interest in maize (Zea mays L.). *Genetics* 190(2):795–811.
83. Pogson GH (1991) Expression of overdominance for specific activity at the phosphoglucomutase-2 locus in the Pacific oyster, *Crassostrea gigas*. *Genetics* 128(1):133–141.
84. Mitchell-Olds T (1995) Interval mapping of viability loci causing heterosis in *Arabidopsis*. *Genetics* 140(3):1105–1109.
85. Graham GI, Wolff DW, Stuber CW (1997) Characterization of a yield quantitative trait locus on chromosome five of maize by fine mapping. *Crop Sci* 37(5):1601–1610.
86. Steinmetz LM, et al. (2002) Dissecting the architecture of a quantitative trait locus in yeast. *Nature* 416(6878):326–330.
87. Stinchcombe JR, et al. (2004) A latitudinal cline in flowering time in *Arabidopsis thaliana* modulated by the flowering time gene *FRIGIDA*. *Proc Natl Acad Sci USA* 101(13):4712–4717.
88. Hancock AM, et al. (2011) Adaptation to climate across the *Arabidopsis thaliana* genome. *Science* 334(6052):83–86.
89. Dittmar EL, Oakley CG, Ågren J, Schemske DW (2014) Flowering time QTL in natural populations of *Arabidopsis thaliana* and implications for their adaptive value. *Mol Ecol* 23(17):4291–4303.
90. Korves TM, et al. (2007) Fitness effects associated with the major flowering time gene FRIGIDA in *Arabidopsis thaliana* in the field. *Am Nat* 169(5):E141–E157.
91. Fournier-Level A, et al. (2013) Paths to selection on life history loci in different natural environments across the native range of *Arabidopsis thaliana*. *Mol Ecol* 22(13):3552–3566.
92. Brachi B, et al. (2010) Linkage and association mapping of *Arabidopsis thaliana* flowering time in nature. *PLoS Genet* 6(5):e1000940.
93. Moore S, Lukens L (2011) An evaluation of *Arabidopsis thaliana* hybrid traits and their genetic control. *G3 (Bethesda)* 1(7):571–579.
94. Lisec J, et al. (2009) Identification of heterotic metabolite QTL in *Arabidopsis thaliana* RIL and IL populations. *Plant J* 59(5):777–788.
95. Meyer RC, et al. (2010) QTL analysis of early stage heterosis for biomass in *Arabidopsis*. *Theor Appl Genet* 120(2):227–237.
96. Oakley CG, Ågren J, Schemske DW (2015) Heterosis and outbreeding depression in crosses between natural populations of *Arabidopsis thaliana*. *Heredity (Edinb)* 115(1):73–82.
97. Schneider CA, Rasband WS, Eliceiri KW (2012) NIH Image to ImageJ: 25 years of image analysis. *Nat Methods* 9(7):671–675.
98. Xiang B, Li BL (2003) Best linear unbiased prediction of clonal breeding values and genetic values from full-sib mating designs. *Can J Res* 33(10):2036–2043.
99. Bates D, Maechler M, Bolker B, Walker S (2015) Fitting linear mixed-effects models using lme4. *J Stat Softw* 67(1):1–48.
100. Devlin B, Roeder K (1999) Genomic control for association studies. *Biometrics* 55(4):997–1004.
101. Chang CC, et al. (2015) Second-generation PLINK: Rising to the challenge of larger and richer datasets. *Gigascience* 4:7.
102. Szkiba D, Kapun M, von Haeseler A, Gallach M (2014) SNP2GO: Functional analysis of genome-wide association studies. *Genetics* 197(1):285–289.
103. Tamura K, et al. (2011) MEGA5: Molecular evolutionary genetics analysis using maximum likelihood, evolutionary distance, and maximum parsimony methods. *Mol Biol Evol* 28(10):2731–2739.