

RESEARCH ARTICLE

# Genetic architecture and selective sweeps after polygenic adaptation to distant trait optima

Markus G. Stetter<sup>1\*</sup>, Kevin Thornton<sup>2</sup>, Jeffrey Ross-Ibarra<sup>1,3\*</sup>

**1** Dept. of Plant Sciences and Center for Population Biology, University of California, Davis, Davis, CA, USA, **2** Dept. of Ecology and Evolutionary Biology, University of California, Irvine, Irvine, CA, USA, **3** Genome Center, University of California, Davis, Davis, CA, USA

\* [mgstetter@gmail.com](mailto:mgstetter@gmail.com) (MGS); [rossibarra@ucdavis.edu](mailto:rossibarra@ucdavis.edu) (JRI)



 OPEN ACCESS

**Citation:** Stetter MG, Thornton K, Ross-Ibarra J (2018) Genetic architecture and selective sweeps after polygenic adaptation to distant trait optima. *PLoS Genet* 14(11): e1007794. <https://doi.org/10.1371/journal.pgen.1007794>

**Editor:** Ryan D Hernandez, UCSF, UNITED STATES

**Received:** June 4, 2018

**Accepted:** October 26, 2018

**Published:** November 19, 2018

**Copyright:** © 2018 Stetter et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** All scripts and code to reproduce the simulations and figures is available at <https://dx.doi.org/10.6084/m9.figshare.6179219>. A detailed interactive graphical analysis of summary statistics is available at <https://mgstetter.shinyapps.io/quantgensimAPP/>.

**Funding:** We acknowledge the financial support of NSF Plant Genome award IOS-1238014, USDA Hatch project CA-D-PLS-2066-H 548 to JRI and the Deutsche Forschungsgemeinschaft (DFG) grant STE 2654/1-1 to MGS. The funders had no role in study design, data collection and analysis,

## Abstract

Understanding the genetic basis of phenotypic adaptation to changing environments is an essential goal of population and quantitative genetics. While technological advances now allow interrogation of genome-wide genotyping data in large panels, our understanding of the process of polygenic adaptation is still limited. To address this limitation, we use extensive forward-time simulation to explore the impacts of variation in demography, trait genetics, and selection on the rate and mode of adaptation and the resulting genetic architecture. We simulate a population adapting to an optimum shift, modeling sequence variation for 20 QTL for each of 12 different demographies for 100 different traits varying in the effect size distribution of new mutations, the strength of stabilizing selection, and the contribution of the genomic background. We then use random forest regression approaches to learn the relative importance of input parameters in determining a number of aspects of the process of adaptation, including the speed of adaptation, the relative frequency of hard sweeps and sweeps from standing variation, or the final genetic architecture of the trait. We find that selective sweeps occur even for traits under relatively weak selection and where the genetic background explains most of the variation. Though most sweeps occur from variation segregating in the ancestral population, new mutations can be important for traits under strong stabilizing selection that undergo a large optimum shift. We also show that population bottlenecks and expansion impact overall genetic variation as well as the relative importance of sweeps from standing variation and the speed with which adaptation can occur. We then compare our results to two traits under selection during maize domestication, showing that our simulations qualitatively recapitulate differences between them. Overall, our results underscore the complex population genetics of individual loci in even relatively simple quantitative trait models, but provide a glimpse into the factors that drive this complexity and the potential of these approaches for understanding polygenic adaptation.

decision to publish, or preparation of the manuscript.

**Competing interests:** The authors have declared that no competing interests exist.

## Author summary

Many traits are controlled by a large number of genes, and environmental changes can lead to shifts in trait optima. How populations adapt to these shifts depends on a number of parameters including the genetic basis of the trait as well as population demography. We simulate a number of trait architectures and population histories to study the genetics of adaptation to distant trait optima. We find that selective sweeps occur even in traits under relatively weak selection and our machine learning analyses find that demography and the effect sizes of mutations have the largest influence on genetic variation after adaptation. Maize domestication is a well suited model for trait adaptation accompanied by demographic changes. We show how two example traits under a maize specific demography adapt to a distant optimum and demonstrate that polygenic adaptation is a well suited model for crop domestication even for traits with major effect loci.

## Introduction

Understanding molecular adaptation is essential for the study of evolutionary processes, genetic diseases, and plant and animal breeding. The process of adaptation is often divided into three separate modes: hard selective sweeps, soft selective sweeps and polygenic adaptation [1]. In recent decades many empirical population genetic analysis have focused on hard selective sweeps because these leave a distinct molecular signature that can be readily detected in genomic data. Hard sweeps result from the reduction of genetic diversity at neutral sites linked to a new beneficial mutation that rapidly fixes [2]. In recent years, other forms of selection that play an important role in evolution and adaptation have begun to receive increased attention, although these are more difficult to detect in empirical data. For instance, sweeps from selection on standing genetic variation leave a less distinct pattern on diversity than hard selective sweeps because the beneficial variant has had more time to recombine onto multiple genetic backgrounds [3, 4]. In addition to processes involving sweeps at individual loci, polygenic adaptation—in which selection acts on a quantitative trait with complex genetic architecture—is frequently regarded as a third mode of adaptation and can lead to rapid phenotypic change via relatively minor shifts in allele frequencies [5].

Although well-studied traits such as human height [6], coat color in mice [7] and grain yield in crops [8] follow patterns consistent with a polygenic pattern, the dynamics and genetic architecture of polygenic adaptation are not well understood. Polygenic adaptation has only gained importance in empirical population genetics relatively recently, but the field of quantitative genetics is based on the idea that traits are controlled by large numbers of loci [9]. Population genetics and quantitative genetics diverged with the appearance of the first molecular data allowing empirical evaluation of single locus population genetic models, while the analysis of effects of single loci in quantitative genetics has long been limited by the large number of phenotyped individuals needed [10]. The increasing availability of high density SNP sets and whole genome sequencing for tens of thousands of individuals, however, is now providing the opportunity to test both population and quantitative evolutionary genetic hypotheses in empirical data [e.g. 11].

Many polygenic traits are thought to evolve under stabilizing selection, in which selection acts against extreme deviations from an optimum trait value [12, 13, 14]. Under such a model, an individual's fitness is given by its phenotypic distance from the trait optimum and the strength of stabilizing selection. Within this framework, recent attention has focused on the dynamics of polygenic adaptation to a new nearby phenotypic optimum [15, 16, 17, 18, 19, 9].

In this scenario, genetic variance in the population decreases when most effect sizes are small, because many sites fix. In contrast, when most mutations have large effect sizes, the genetic variance increases because large effect loci increase in frequency but do not fix [15, 19]. In addition to allele frequency changes, theoretical quantitative genetic analyses have revealed that selective sweeps are prevalent during polygenic adaptation [20, 17]. These studies have developed important theoretical background for the understanding of polygenic adaptation and have documented the dynamics of a small number of loci during the course of adaptation. Each of these studies shows in detail how a small number of parameters influences adaptation, but the complex interplay of mutation, selection, and demography across a large parameter space has not yet been explored. For example, population growth has been shown to influence the contribution of low frequency alleles to trait variance [21], but the interaction of demography with parameters such as the distribution of effect sizes of new mutations needs further investigation.

Here, we take a simulation approach to study a population adapting to an optimum shift, modeling sequence variation for 20 QTL for each of 12 different demographic models for 100 different traits with varying effect size distribution of new mutations, strength of stabilizing selection, and the contribution of the genomic background beyond the simulated QTL. After detailed analysis of a single scenario, we use machine learning to extract parameter importance for the input parameters. Our results illustrate that selective sweeps are common under most scenarios, even for mutations of relatively minor effect. We employ machine learning on genetic architecture matrices and find that demography and the effect size of new mutations have the largest influence on present day genetic architecture. After identifying general parameter importance, we use maize domestication as an example and investigate two diverging traits in a population that underwent a population bottleneck and exponential growth [22], showing how these traits adapt to the changing optimum and comparing our findings to archaeological and genetic data [23, 24].

## Results

We first simulated adaptive and stabilizing selection on a single quantitative trait in a randomly mating diploid population. After a burn-in to equilibrium, we simulated an instantaneous shift in the optimal trait value from 0 to 10, corresponding to 89.6 z-scores of  $V_{G_0}$ . The population underwent truncation selection until reaching the new optimum, at which time stabilizing selection resumed. We assumed an additive model with no epistasis, and simulated 20 unlinked QTL as well as a genomic “background” over a range of parameters describing population demography and the trait, including the effect size of new mutations, strength of stabilizing selection, distance to the new optimum, effects of genomic background, and bottleneck severity and population expansion (Table 1).

**Table 1. Parameters and variables.**

Variable	Description
$N_{anc}$	Population size at equilibrium
$N_{final}$	Population size after $0.1 \times N_{anc}$ generations
$N_{bottleneck}$	Population size during bottleneck
$\psi$	Proportion of phenotype due to genetic background outside of QTL
$\sigma_m$	Standard deviation of effect sizes of new mutations
$V_S$	Strength of stabilizing selection
$V_{G_0}$	Genetic variance at equilibrium

<https://doi.org/10.1371/journal.pgen.1007794.t001>

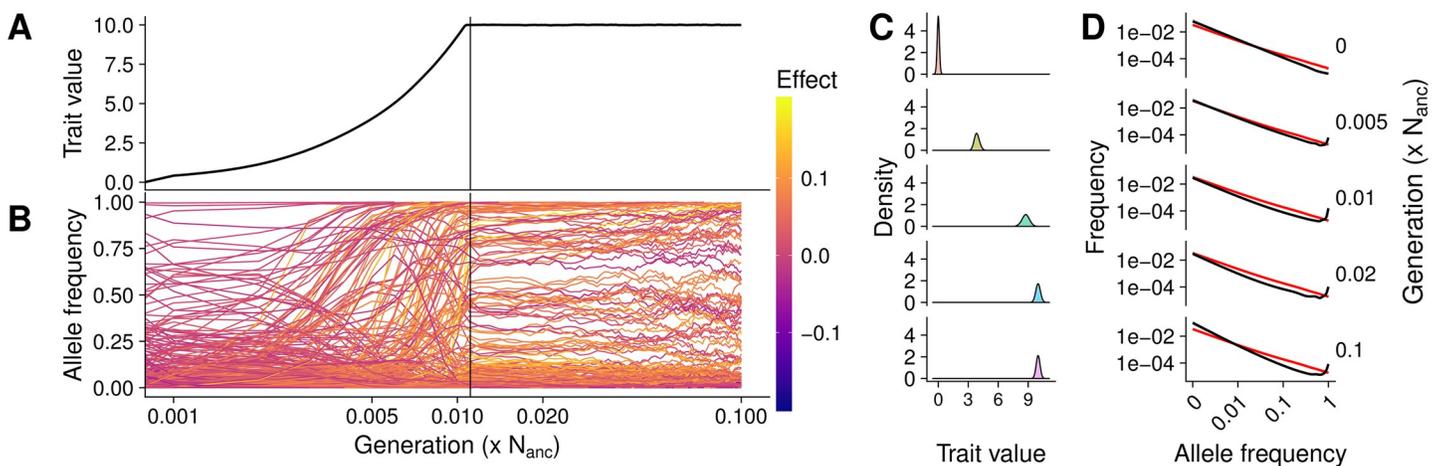
### Single simulation results

The adaptation of a quantitative trait to a sudden environmental change involves allele frequency shifts at many sites, some of which result in selective sweeps. To build intuition around basic patterns seen in these simulations as a population adapts to a new optimum, we first describe results of a single simulation with constant population size, intermediate effect sizes of new mutations ( $\sigma_m = 0.05$ ), strong stabilizing selection ( $V_S = 1$ ), and no phenotypic effect of the genomic background ( $\psi = 0$ ). We present how such a population adapts to the new optimum and how allele frequencies and effect sizes change during this process (Fig 1).

The population mean trait value increased linearly ( $\log_{10}$  scaled x-axis in Fig 1A) until shortly before the new optimum was reached within  $0.011$  ( $sd = 0.0004$ )  $N_{anc}$  generations (Fig 1A and 1C). As the population mean approached the optimum the rate of change decelerated, presumably because some individuals now had phenotype values above the optimum such that alleles which contribute positively to the trait are no longer uniformly beneficial to fitness. The trait variance increased after the optimum shift and during the adaptation process. Though it declined once the new optimum was reached, it did not return to the equilibrium variance by the end of the simulation (Fig 1C). This increase in variance is generated by the increase in allele frequency of formerly rare, large positive effect alleles.

Following individual mutations shows that, at the onset of the optimum shift (generation 0) alleles with negative effect sizes rapidly decline in frequency unless they were already near fixation (Fig 1B). Alleles with positive effects, on the other hand, increase quickly in frequency and fix. Once the new optimum is reached, frequencies of both positive and negative alleles changed slowly, but the number of small effect alleles increased. This shows how a population can adapt to a sudden environmental change by an increase of beneficial alleles and decrease of disadvantageous alleles in a relatively short time.

Looking at the change in allele frequencies of all mutations helps to understand what drove the adaptation process in the population (Fig 1D). At equilibrium, variants with larger effects



**Fig 1. Population dynamics of a single parameter set.** A) Trait evolution after an optimum shift and B) allele frequency dynamics during adaptation from a single replicate. The vertical line shows when the new trait optimum was reached and line colors denote effect sizes. Time is shown on a log scale. C) The phenotypic distribution and D) site frequency spectra of segregating mutations (black) and neutral expectation (red) from 100 independent replicates. Panels show different generations including equilibrium prior to adaptation (0), during adaptation (0.005), just before the new optimum is reached (0.01), after the new optimum has been reached (0.02), and the final generation (0.1). All results are from a simulated population with constant population size,  $\sigma_m = 0.05$ , and  $V_S = 1$ .

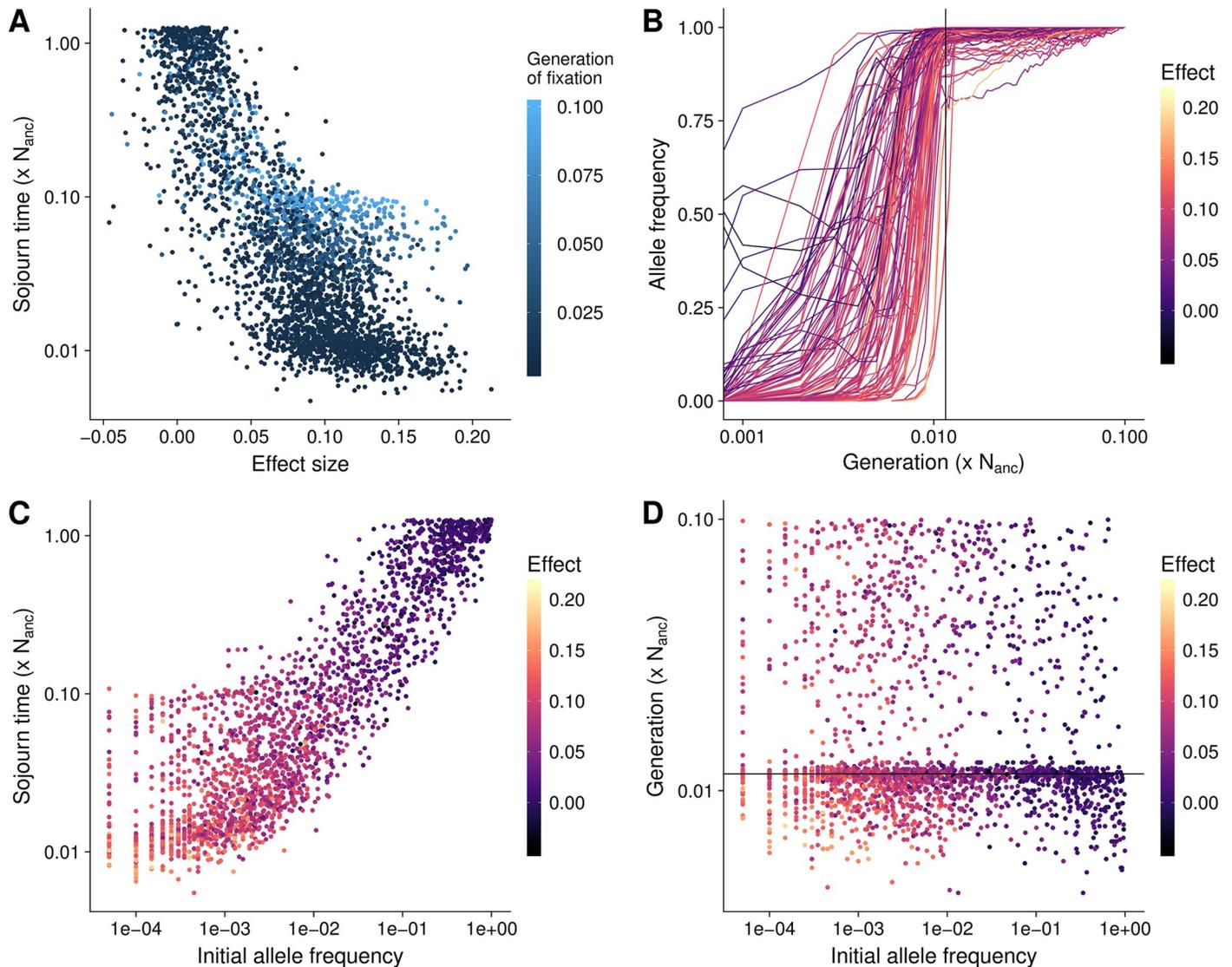
<https://doi.org/10.1371/journal.pgen.1007794.g001>

are selected against, leading to an excess of rare variants compared to neutral expectations. The site frequency spectrum (SFS) then changed quickly after the optimum shift as selection fixed positive mutations. Directly before the new optimum was reached ( $0.01N_{\text{anc}}$ ), 11% of mutations were at very high frequencies ( $> 0.5$ ) while after reaching the new optimum ( $0.02N_{\text{anc}}$ ) only 8% of mutations were at such high frequencies and the number of high frequency segregating sites further declined in consecutive generations. Under stabilizing selection, extreme values are again selected against and alleles that have risen to intermediate frequency during adaptation fix or get lost. By  $0.1 \times N_{\text{anc}}$  generations the SFS again reflected an excess of rare alleles, but also an excess of high frequency derived alleles. The observed high frequency derived alleles ((Fig 1D) bottom) represent in fact the other allele, which is at low frequency. These mutations increased in frequency during adaptation, but both alleles have the same fitness effect after the equilibrium has been reached and the mutation does consequently not decrease in frequency.

When a selected mutation increases in frequency quickly, it often reduces diversity in adjacent genomic regions, leading to a pattern commonly referred to as a selective sweep. While we cannot assess diversity at linked neutral sites in our model, we can nonetheless identify likely selective sweeps by comparing the sojourn time of individual alleles to that of a neutral allele experiencing equivalent demographic processes (see [Methods](#)). Following these criteria, 72% of all fixations in this simulation were selective sweeps. Of these, 73% were sweeps from standing variation. While there was an overall negative correlation between the time a site was segregating in the population and its effect size on the trait, there were a number of mutations that fixed later than expected given their effect size (Fig 2A).

Observing the frequency trajectories of sites that fixed after the new optimum had been reached shows that the speed of frequency change slowed down substantially, but these alleles eventually reached fixation. When the new optimum has been reached, any increase or decrease in frequency of large effect mutations takes the population away from the trait optimum and is selected against. The remaining change in frequency is mostly stochastic and results from minor fluctuations in the trait mean due to frequency changes at other sites [19]. However, because stabilizing selection acts against stochastic variation in allele frequencies that move the population away from the optimum, the time to fixation or loss for an allele is still faster than neutrality in a manner that has sometimes been deemed similar to underdominance [25]. Some mutations with negative effects that decreased in frequency under truncation selection after the optimum shift can then increase in frequency again once the new optimum is reached and stabilizing selection takes over (Fig 2B). Such mutations provide a good example of selection on a quantitative trait, which results in selection coefficients that can vary in sign or magnitude depending on the total phenotypic value of the individual in which they occur, its distance to the optimum, and the details of when and what kind of selection occurred.

In our simulations, fixations from standing variation fixed either fast, because they were present at high frequency at the onset of directional selection, or due to their large effect on the trait. However, there was no correlation between the initial allele frequency and the generation in which the mutation fixed (Fig 2C and 2D). Large effect mutations segregated at low frequency in the equilibrium population, while small effect sites were already at higher frequencies, explaining why large effect and small effect mutations fixed at similar generations, despite the difference in speed of allele frequency shift. Negative and effectively neutral mutations may also fix together with large effect positive mutations presumably due to the effects of genetic hitchhiking (Fig 2).



**Fig 2. Selective sweeps.** A) Speed of fixation of selective sweep mutations. B) Dynamics of fixations that occur after the new optimum was reached. C) Speed of fixation of sweeps from standing variation compared to their initial frequency. D) The generation at which sweeps from standing variation fix. All results are from a simulated population with constant population size,  $\sigma_m = 0.05$ , and  $V_S = 1$ , and time is shown on a log scale.

<https://doi.org/10.1371/journal.pgen.1007794.g002>

### Complex genetic architectures with demographic changes

The detailed analysis of a single population adapting to a sudden environmental change helps to build intuition on the dynamics of a specific set parameters, but is far from the complexities of quantitative trait evolution in natural populations. For example, most populations have experienced some form of fluctuation in population size, and traits differ both in the strength of stabilizing selection as well as in their genetic architecture—the frequency and effect size of mutations that cause variation in the phenotype. To understand the effect of these and other variables, we simulated 1,200 different combinations of parameter sets to examine the contribution of the strength of stabilizing selection, the effect size of new mutations, population demography, and differences in genetic background on variation and adaptation of the focal trait.

The combination of  $V_S$  and  $\sigma_m$  led to different genetic variances at equilibrium ranging from 0.004 to 0.751, leading to a distance of the new trait optimum between 11.5 and 158.2 z-scores (S3 Fig). We calculate  $V_G$  in every generation during the burn-in and compared it to the expected genetic variance under the House of Cards (HoC) or stochastic HoC [26, 27] approximations [12]. The majority of simulations are within the regime of HoC, though the approximation underestimated  $V_G$  for  $\sigma_m = 0.9$  and  $V_S = 1$  and overestimated  $V_G$  for large  $V_S$  and small  $\sigma_m$ ; all simulations closely matched expectations under the stochastic HoC approximation. All burn-ins had a mean fitness close to one at equilibrium after  $10N$  generations and the mean  $V_G$  was constant (S4 and S3 Figs).

To understand the factors driving variation in particular aspects of the data, we employed a random forest machine learning model (see Methods) to retrieve parameter importance.

### Speed of polygenic adaptation

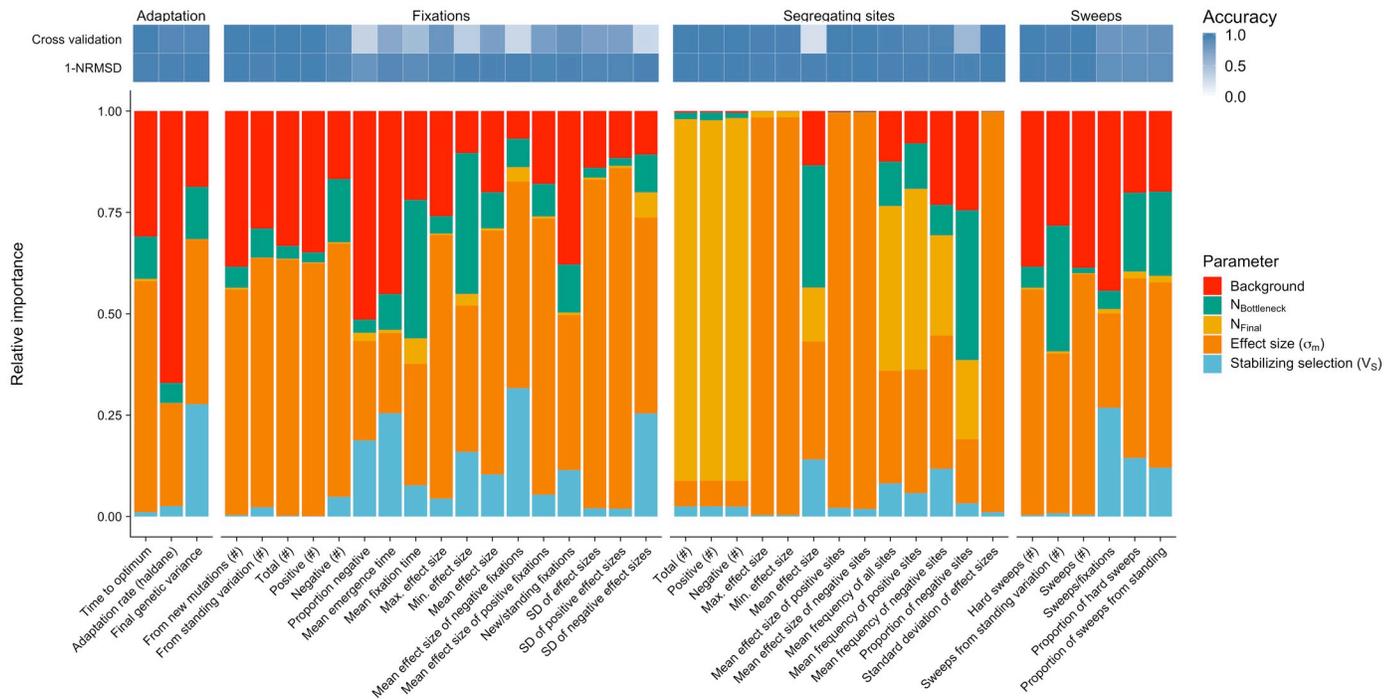
An important factor for the survival of a population exposed to changing environments is how fast it can adapt to new conditions. Our simulated populations varied widely in the time required to reach the new optimum, from 0.001 to 0.099  $N_{\text{anc}}$  generations. A total of 732 of the 120,000 simulations did not reach the new trait optimum within the simulated time of  $0.1 \times N_{\text{anc}}$  generations, but all parameter combinations had at least 8 (of 100) replicates reaching the new optimum. In general, simulations that did not reach the new optimum were those with a strong bottleneck (reduction to 1% or 5% of  $N_{\text{anc}}$ ). In particular, more than 70% of all simulations with the smallest  $\sigma_m$  (0.01), no genetic background, 1% bottleneck, and a final size of  $N_{\text{anc}}$  did not reach the new optimum, regardless of their strength of stabilizing selection ( $V_S$ ).

All three adaptation-related summary statistics (time to new optimum, adaptation rate, and  $V_G$  in the final generation) were well predicted, with cross-validation accuracy over 90% (Fig 3 top). Overall, the parameter contributing most to this variation is  $\sigma_m$ , with a relative importance of > 50% (Fig 4). This was followed closely by the proportion of the trait explained by genetic background ( $\psi$ ) at 31%, while demography and  $V_S$  were of relatively minor importance (Fig 3 and S5 Fig). We find that the rate of phenotypic adaptation was highest for populations with small  $\sigma_m$  and large  $\psi$ , and these two factors explained the majority of the observed variation (Fig 4). The initial genetic variance, a combination of  $V_S$  and  $\sigma_m$ , was the best predictor for the genetic variance in the final generation, but the strength of the bottleneck and  $\psi$  had a relative importance of 11% and 17%, respectively (S5 Fig). The genetic variance in the final generation increased with increasing  $\sigma_m$ , though it plateaued at the largest  $\sigma_m$  simulated (Fig 4). Of the 1,200 parameter combinations, 45 had a mean  $V_G$  as high as the initial  $V_G$  ( $\pm 5\%$ ), 410 had higher  $V_G$  and 745 had lower  $V_G$  than the initial equilibrium population.

### Segregating sites after polygenic adaptation

We further investigated segregating sites in the final generation, which correspond to a modern population that has experienced an optimum shift in the past. Cross validation prediction accuracies were for most summary statistics very high ( $<0.9$ ). The mean effect size of segregating sites was predicted with less accuracy, however, as all values are concentrated around zero leading to low  $R^2$  values in the CV. The NRMSE, shows that the accuracy for mean effect size of segregating sites was high and that the validation data could be predicted, which allowed to infer parameter importance even with lower CV accuracy.

While absolute numbers mostly depended on the final population size, other statistics showed more distinct patterns. Allele frequencies of both negative and positive sites were strongly influenced by the demography of the population. The proportion of negative sites



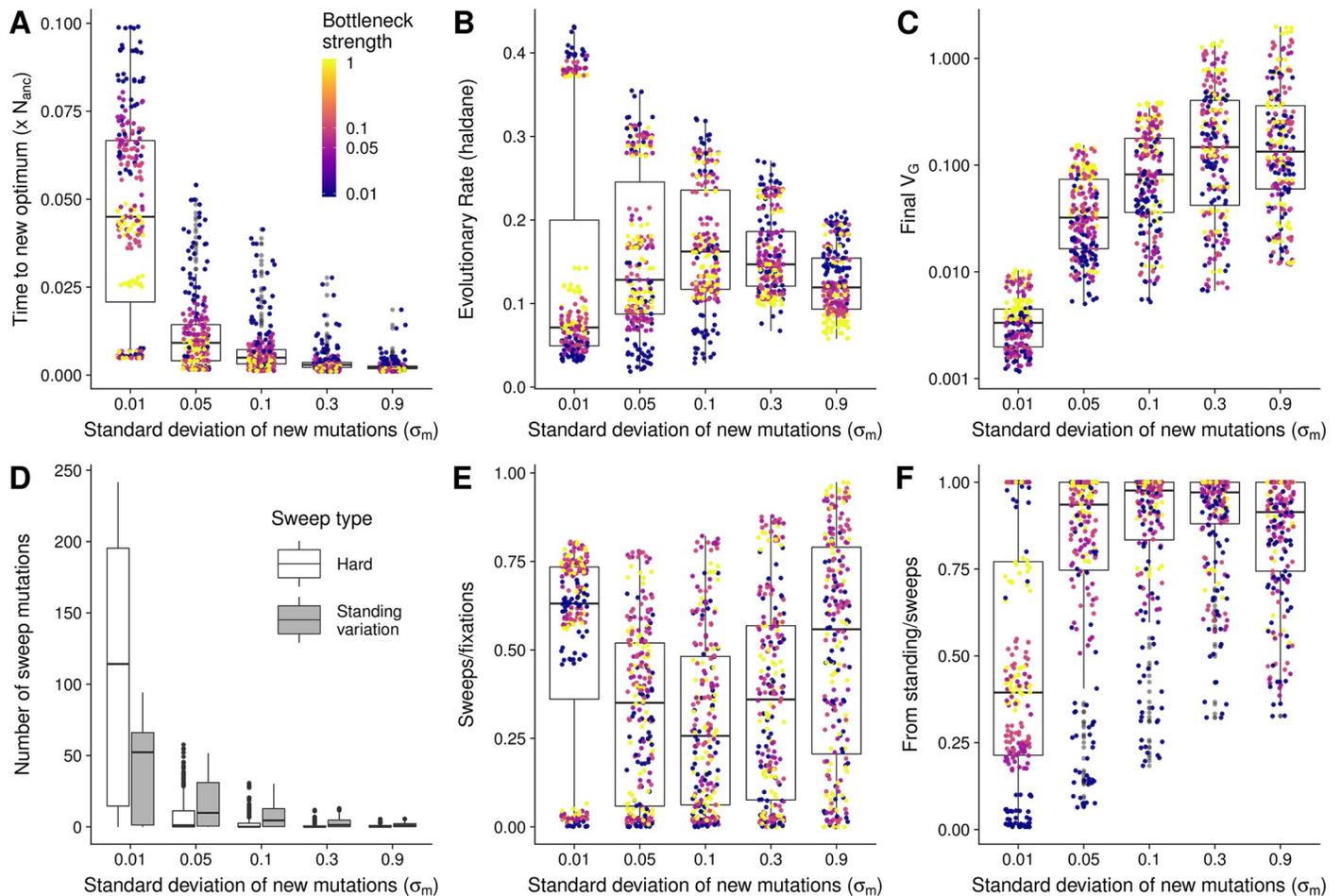
**Fig 3. Relative parameter importance.** Relative parameter importance inferred for four parameter categories. 1) Adaptation: parameters describing adaptation speed and potential for future adaptation, 2) Fixations: summary statistics for mutations that were fixed during trait adaptation, and 3) Segregating sites: descriptors of alleles polymorphic in the final generation of the simulations. Top rows indicate prediction accuracy as calculated by 10-fold cross validation and NRMSE. Each bar is the result of an independent random forest learning and each color represents the relative importance of the simulation input parameters (see [Methods](#) and [S1 Table](#) for summary statistics).

<https://doi.org/10.1371/journal.pgen.1007794.g003>

segregating in the population was also most strongly influenced by the strength of the bottleneck (Fig 3), but when  $V_{G_0}$  (S3 Fig) was used to train the model instead,  $V_{G_0}$  explained most of the variation (S5 Fig). As  $V_{G_0}$  is the result of the combination of  $V_S$  and  $\sigma_m$  during the burn-in, there is a strong interaction effect between  $V_S$  and  $\sigma_m$ , which is partitioned when using  $V_S$  as feature in the random forest.

### Fixations and selective sweeps

Mutations in a population can rise in frequency and fix due to demographic events and stochastic sampling or as a result of selection. The sudden change in trait optimum in our model imposed strong selection on sites with a positive effect, while mutations with negative effect values were deleterious until the new optimum was reached. Different parameter combinations led to strongly varying numbers and patterns of fixations in our simulations. The effect size of new mutations ( $\sigma_m$ ) and  $\psi$  had the strongest influence on the absolute number of fixations and the effect size of mutations that fixed (Figs 4, 3 and S5 Fig). Variation in the mean effect size of fixations depended mostly on  $\sigma_m$ , though  $V_S$  also contributed substantially for negative fixations. Consistent with fixations being driven primarily by selection, the effect size of positive mutations that fixed was an order of magnitude larger than that of negative fixations (S6 Fig). Comparing results within each set of simulations with identical  $\sigma_m$  shows that stochastic sampling due to  $N_{bneck}$  played an important role in determining the number of fixations even if the relative importance of  $N_{bneck}$  among all parameters was only 3% (S6A Fig and Fig 3).



**Fig 4. Summary of trait adaptation and selective sweeps.** A) Time to reach new trait optimum B) Rate of change in phenotype C) Genetic variance after  $0.1 \times N_{anc}$  generations. D) Total number of selective sweeps, separated by type of sweeps. E) Proportion of sweeps compared to all fixations F) Proportion of sweeps from standing variation. Boxes are split by major parameter importance as identified by our random forest model. Points in A-C and E-F show the values of each of 1,200 parameter sets and are colored according to bottleneck size (Darker color indicate stronger bottleneck, see legend in A). Interactive plots are available at <https://mgstetter.shinyapps.io/quantgensimAPP/>.

<https://doi.org/10.1371/journal.pgen.1007794.g004>

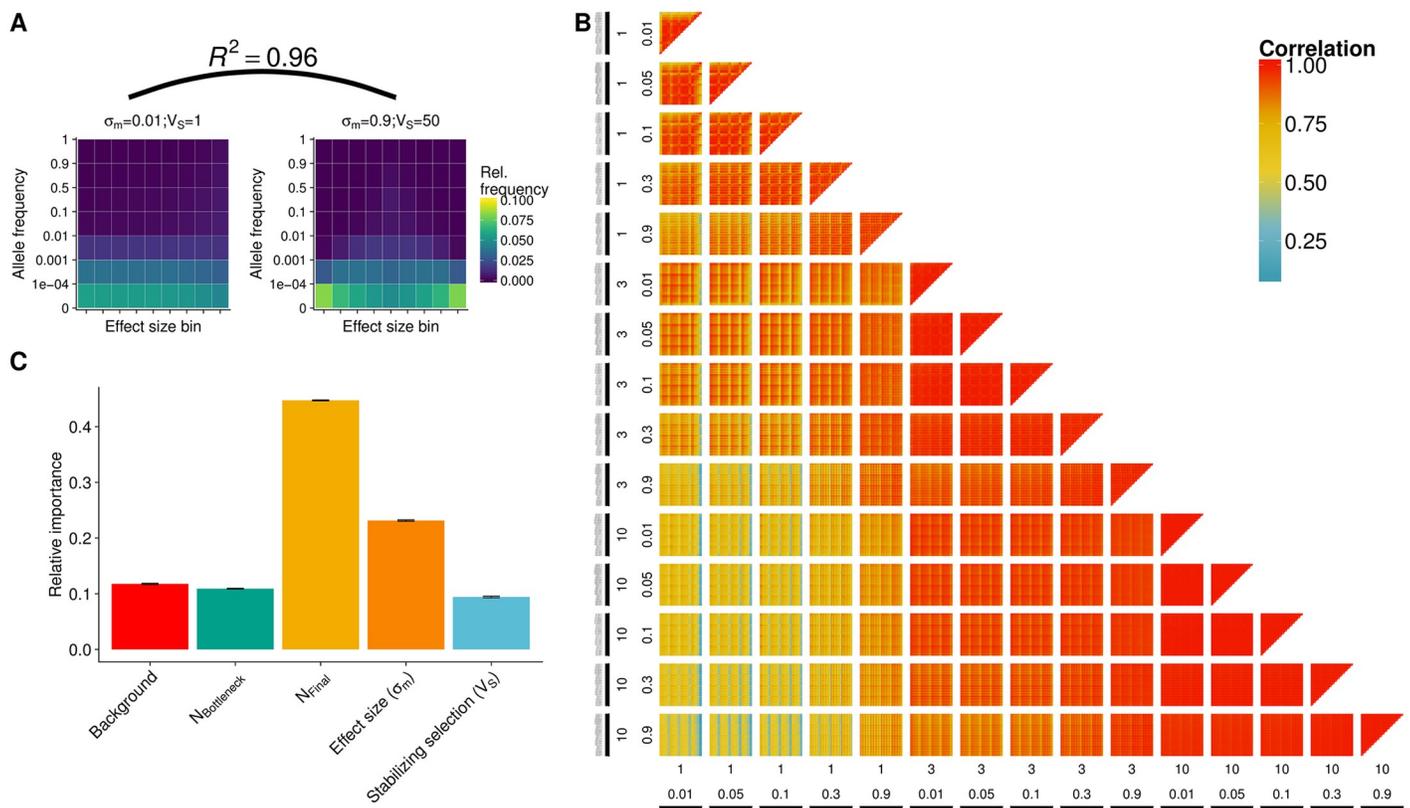
Not all fixations are due to positive selection, however, and even those that are due to selection would not necessarily reduce linked diversity sufficiently to be detected as a selective sweep. To differentiate between neutral and strongly selected fixations, we compared the fixation time of sites that fixed after the shift in trait optimum to single-locus neutral simulations with identical demography (see [Methods](#)). Consistent with the higher total number of fixations exhibited, populations with smaller  $\sigma_m$  also showed a higher number of sweeps. While the maximum number of sweeps was almost 300 (for  $\sigma_m = 0.01$ ,  $\psi = 0$ ,  $V_S = 1$ , and a bottleneck), 13 parameter sets did not lead to any sweeps within the simulated time, all with  $\sigma_m \geq 0.3$ ,  $\psi = 0.95$  and  $V_S \geq 5$ . The proportion of sweeps to fixations ranged from 0 to 99% but was highly variable and revealed strong interactions between  $\sigma_m$ ,  $\psi$  and  $V_S$  (Fig 4). Larger  $\psi$  led to a low proportion of sweeps to fixations when  $V_S$  and  $\sigma_m$  were small, but for large values of  $V_S$  and  $\sigma_m$  almost all fixations were sweeps, scaling with decreasing  $\psi$  (<https://mgstetter.shinyapps.io/quantgensimAPP/>). The proportion of sweeps from standing variation was also highly variable, but differentiated more strongly by demography within each group of  $\sigma_m$  than the total proportion of sweeps (Fig 4E). Population bottlenecks were the second most important

parameter for the type of selective sweep observed, while either  $\sigma_m$  or  $V_{G_0}$  were the most important parameters (Fig 3 and S5 Fig).

### Genetic architecture after adaptation

The genetic architecture of phenotypic traits that we observe in populations today was shaped by demographic history and past selection. We evaluated the genetic architecture in the final generation of all 1,200 populations with their diverse range of histories by comparing the combined allele frequency—effect size matrices (see Methods). These matrices were used as input for our random forest model to understand the contributions of input parameters to variation in genetic architecture.

The extracted parameter importance showed that the variation in the genetic architecture depended most strongly on  $N_{\text{final}}$  and  $\sigma_m$ , but each of the other three parameters contributed at least 9% of the variation (Fig 5). The strong interaction between parameters becomes apparent in Fig 5, where the fine structure beyond the major 2 parameters ( $\sigma_m$  and  $N_{\text{final}}$ ) can be seen on all levels of combinations. Among simulations with large  $\sigma_m$  and large  $N_{\text{final}}$ , however, all correlations are close to 1 and it is therefore not possible to easily distinguish parameter sets based on their genetic architecture (individual genetic architecture plots for each parameter combination are available at <https://mgstetter.shinyapps.io/quantgensimAPP>).



**Fig 5. Genetic architecture in final population.** **A**) Genetic architecture matrices for two parameter combinations (maize models, see Methods) differing in effect size of new mutations and strength of stabilizing selection. Effect size bins are centered around zero with negative effect size quantiles on the left and positive quantiles on the right of the central bin. Shown is the correlation coefficient between the genetic architectures. **B**) Pairwise correlation of genetic architecture of all comparisons of 1,200 parameter combinations. Subplots display the combination of final population size (log: 1, 3, 10) and effect size distribution ( $\sigma_m$ , 0.01, 0.05, 0.1, 0.3, 0.9) of incoming mutations. Each pixel displays a pairwise comparison between two of the 1,200 scenarios. **C**) Relative parameter importance for genetic architecture prediction.

<https://doi.org/10.1371/journal.pgen.1007794.g005>

## Maize domestication traits

After evaluating a wide parameter space using our machine learning models, we then investigated in more detail two parameter sets that resemble diverging traits during maize domestication. Using simulations with demographic models similar to that inferred for maize [a bottleneck of  $0.05 \times N_{\text{anc}}$  followed by exponential growth to  $10 \times N_{\text{anc}}$ , 22], we selected one trait with strong stabilizing selection and small effect mutations (Trait 1;  $\sigma_m = 0.01$  and  $V_s = 1$ ) and one trait with weak stabilizing selection and large effect mutations (Trait 2;  $\sigma_m = 0.9$  and  $V_s = 50$ ).

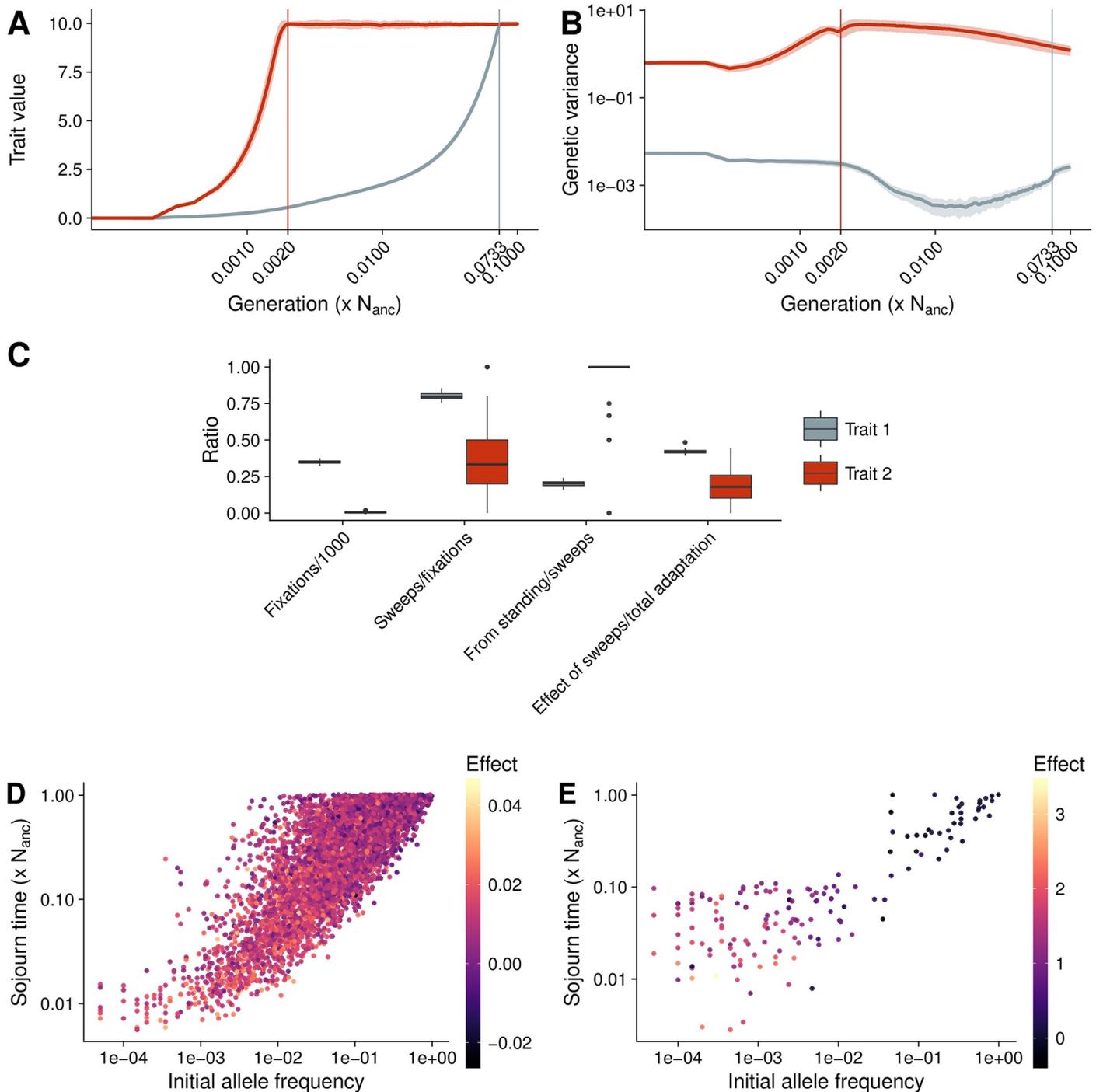
The two traits showed notably different patterns of adaptation (Fig 6, x-axis on  $\log_{10}$  scale). Trait 1 increased almost linearly for  $0.0733 \times N_{\text{anc}}$  generations before asymptotically arriving at the new optimum. The genetic variance for this trait declined for the first  $0.0169 \times N_{\text{anc}}$  generations before it slowly increased, but did not reach the equilibrium value within the  $0.1 \times N_{\text{anc}}$  generations simulated. Trait 2, on the other hand, adapted rapidly, reaching the optimum in only  $0.002 \times N_{\text{anc}}$  generations. The genetic variance for Trait 2 increased during adaptation to a value higher than  $V_{G_0}$ , then decreased after the optimum was reached but remained higher than  $V_{G_0}$  (Fig 6A and 6B). The number of fixations was 100 times higher for Trait 1 than for Trait 2; the ratio of sweeps per fixation was also higher, and most sweeps in Trait 1 were hard (Fig 6C). Though on average Trait 2 exhibited fewer than 2 sweeps per simulation, 94% of these were from standing variation. Neither trait showed a strong correlation between mutation effect size and when fixation occurred, suggesting that the domestication bottleneck was not the primary driver of fixation (S7 Fig). The sojourn time for sweeps from standing variation was correlated with the initial allele frequency, but also with its effect size. Large effect positive mutations had a low initial frequency but fixed quickly, while negative alleles fixed slowly despite their high initial frequency, similar to the initial trait described above (Fig 2). This observation held particularly true for Trait 2, where only few small or negative effects fixed quickly (Fig 6D and 6E). The overall contribution of all sweeps to phenotypic change was also different between the two traits: the summed effect size of all sweeps represents 45% of the adaptation in Trait 1, but only 18% for Trait 2.

Fig 5A shows the difference in genetic architecture between the two traits. While the adaptation of Trait 1 led to an equal distribution of effect sizes at low frequencies, Trait 2 had a larger proportion of both very low frequency mutations from the extreme tail of the distribution and small effect mutations at higher frequencies. Despite these differences the correlation between the genetic architecture matrices was very high (0.96; Fig 5).

## Discussion

### Model choice

We use a combination of two different fitness functions to study the quantitative genetics of adaptation to a sudden change to a new trait optimum far beyond observed trait values for any individual in the equilibrium population. During the stationary phase before the shift and after reaching the new optimum we followed a Gaussian fitness function appropriate for a trait under stabilizing selection [14]. During the optimum shift, however, such a model would be problematic, as only a few individuals in the upper tail of the fitness distribution would have extremely high relative fitness, inducing a strong population bottleneck. Instead, we applied a model of truncation selection, first calculating fitness under the Gaussian fitness function but then assigning a fitness of 1 to the top half of the population and 0 to the bottom half. Such a model is reasonable for sudden shifts in trait optima that do not lead to the extinction of a population, but where higher trait values are unambiguously advantageous and the maximum



**Fig 6. Maize specific adaptation.** A) The evolution of trait value and B) genetic variance during adaptation to a new trait optimum for two traits under maize demography with no genetic background. Time in both figures is shown on a log scale, light shadows show standard deviations from the mean of 100 simulation replicates. Trait 1 (blue) has small effect mutations ( $\sigma_m = 0.01$ ) and strong stabilizing selection ( $V_S = 1$ ). Trait 2 (red) has large effect new mutations ( $\sigma_m = 0.9$ ) and weak stabilizing selection ( $V_S = 50$ ). Vertical lines denote the generation when 99% of the new trait optimum is reached. C) Proportion of selective sweeps. D) Sojourn time of sweeps from standing variation in Trait 1. E) Sojourn time of sweeps from standing variation in Trait 2. Scales in D and E are different due to strong divergence of effect size values.

<https://doi.org/10.1371/journal.pgen.1007794.g006>

population size is limited. In natural populations these factors can be observed when sudden changes in the environment favor a specific phenotype for invasive species [28] or in semi-artificial populations in agroecosystems and during domestication [24]. Truncation selection is also common in evolve and re-sequence experiments [29], crop populations [30] and during strong directional selection in natural populations [31].

In our model simulations we fixed the equilibrium optimum to 0 and the new optimum to 10, but change  $V_S$  and  $\sigma_m$ . As  $V_S$  and  $\sigma_m$  change, the relative distance to the new optimum changed with respect to the initial  $V_G$  ( $V_{G_0}$ ). The wide range of distances simulated resembles observations in nature and experimental populations. For example, in the Illinois long term selection experiment in maize, 105 generations of selection for high oil resulted in a shift of over 40 standard deviations [30], and large trait shifts have also been identified in other experimental and natural populations [32, 33]. Our results should therefore be relevant for a variety of traits that adapt to changing environments.

While our modeling investigated a wide parameter space for a number of key variables, one key aspect we have ignored is interaction among alleles (dominance) or loci (epistasis). Both forms of interaction are widely recognized to be important at the molecular level [34, 35], but the majority of variance for a wide array of quantitative traits seems reasonably well explained under a simple additive model [36, 37], but see [38, 8, 34, 39]. Although we do not include any explicit simulation of interlocus interactions, our quantitative trait model is such that the effect of an allele in any given generation will depend on the genetic background. We predict that epistasis and dominance would absorb some of the effect of  $\sigma_m$  for most statistics and have relatively little influence on demographic parameters. Further efforts should incorporate the effects of dominance and epistasis, especially for understanding phenomena such as heterosis and inbreeding depression, where non-additive effects are likely to play a significant role [40, 41].

## How do organisms adapt to change

Rapidly changing environments, such as those faced by changing climate, impose a threat to populations with narrow genetic variance for important traits [42]. However, the speed and manner in which traits adapt depend on the initial variation and beneficial mutations entering the population once the environment changes. In rapidly changing environments or during new colonization of habitats the time it takes to reach the new optimum is critical, as this might determine whether the population is first to occupy a niche. We looked at two summary statistics—time to optimum and adaptation rate—to compare the adaptive behavior of different traits. The speed to the optimum shows the absolute speed of a population to reach the new optimum, while the adaptation rate is corrected for the genetic variance present. The absolute speed depends most on  $\sigma_m$ , but the adaptation rate is more uniform across  $\sigma_m$  with even higher adaptation rates for small  $\sigma_m$  (Fig 4A and 4B). This shows that with small effect mutations and strong stabilizing selection adaptation is mutation limited, but this is not the case when  $V_S$  is large. These two types of adaptation regimes have previously been described as mutation and environmentally limited adaptation regime [16]. Large adaptation rates are reached with the largest  $\psi$  (0.95) values, because genetic diversity is maintained during the adaptation process. Populations with small  $\sigma_m$  and small  $\psi$  run out of genetic variance, because most positive standing variation fixes and negative mutations get lost. The loss of genetic variance is also apparent when comparing the initial genetic variance to the final genetic variance, which is smaller after adaptation for most populations with  $\sigma_m = 0.01$  (Figs 4C and 6B and S3 Fig). The decrease in  $V_G$  for small effect mutations and the increase from large  $\sigma_m$  is consistent with previous results [15]. The genetic variance after historical adaptation is important in the

face of climate change where recently adapted populations will be forced to further adapt. Populations with a large initial genetic variance and large effects also have larger genetic variance in the final population and are thus better prepared for future adaptation. The severity of population bottlenecks is an additional factor influencing  $V_G$  in the final generation as diversity is removed by genetic drift (Fig 3 and S5 Fig).

Overall, populations with the largest  $V_{G_0}$  and largest  $\sigma_m$  adapt fastest to a new optimum as expected, but we also show the impact of population bottlenecks and the overlap between trait architectures (combinations of  $V_S$  and  $\sigma_m$ ). Different trait architectures can result in similar adaptation speed and genetic variance depending on the population history. This implies that for traits that are highly polygenic, it is of particular importance to prevent population declines in order to maintain the adaptability of populations.

### Selective sweeps during polygenic adaptation

Much of standard population genetic theory assumes mutations have a constant fitness effect  $s$ . This assumption has led to a number of findings about selective sweeps, from the probability of fixation being  $\approx 2s$  [43] to the rule of thumb that mutations with fitness effects  $|2Ns| > 1$  will be fixed or removed by natural selection, while those with smaller effects will drift stochastically as effectively neutral alleles [44]. For quantitative traits, however, the fitness effect of a mutation is conditional on the phenotypic distance of an individual to the trait optimum and the correlation between the trait and fitness [14]. At equilibrium this follows a Gaussian distribution (Eq 1), but during directional selection it will depend on the distance of the population from the trait optimum. The relationship between the phenotypic effect size of a mutation and its fitness effect is strongly positive at the onset of selection, while the slope declines as the population trait mean approaches the new optimum and is even slightly negative once the new optimum has been reached [17, S8 Fig]. This shows that segregating large effect positive mutations are beneficial when the population trait mean is distant from the new optimum, but become disadvantageous once the population mean is close to the new optimum, as on average they will cause individuals to overshoot the optimum.

Most selective sweeps occur during the adaptation process before the new optimum has been reached, but the number of fixations and sweeps is strongly dependent on the demography of the population. A strong population bottleneck leads to more fixations, but most of these are fixed by drift rather than selection, and  $N_{\text{bottleneck}}$  is therefore more important for the number of fixations than the number of selective sweeps (Fig 3 and S5 Fig). Population bottlenecks also decrease the proportion of sweeps from standing variation and favor hard selective sweeps, because the bottleneck removes segregating beneficial alleles (Fig 4).

The overall importance of selective sweeps for different traits depends on the initial genetic architecture: our two example traits show that differences in the number of sweeps do not necessarily reflect their combined effect. Trait 1 exhibited 279 sweeps, these contributed to 42% of the change in trait value, while for Trait 2, only 2 sweeps contributed 18% (Fig 6C). This is consistent with previous results showing that allele frequency shifts of large effect alleles are sufficient to reach the new optimum, but selective sweeps are more important when the new optimum is distant [20, 15]. Our results show even more extreme cases, for example Trait 1 and simulations with  $\sigma_m \leq 0.05$ , in which the population exhausts standing variation and relies almost entirely on new mutations. In this case hard selective sweeps are most common, as new positive mutations provide a strong relative fitness advantage (Figs 4 and 6).

Without linked neutral sites, our ability to identify likely sweep regions requires a few important caveats. First, we use a conservative definition of selective sweeps, including only those alleles fixing faster than 99% of neutral simulations. Less conservative cutoffs should not

strongly influence the general result, as most mutations that sweep fix substantially faster than neutral fixations and only a few more fixations would be defined as selective sweeps. Second, while we identify only sweeps from mutations that arose after the optimum shift as hard sweeps, some sweeps from standing variation would be difficult to distinguish from hard sweeps in genomic data if their frequency at the onset of directional selection was very low [45]. Likewise, not all alleles that fixed faster than 99% of neutral simulations would be detectable as selective sweeps in empirical data, as selection on standing variation has a less pronounced impact on diversity at linked sites [4].

### The effect of genetic background on focal QTL

Allele frequency shifts and selective sweeps in a focal QTL are dependent on the genetic background. Chevin and Hospital [17] showed analytical results for the behavior of a single locus with a polygenic background during the adaptation to a new optimum. In our study, we simulate a more complex case: in addition to a genetic background (see Eq 4), we model 20 QTL each involving numerous loci. Moreover, in our model the QTL and the genetic background are not independent, because the QTL in the parents contribute to their trait value but can themselves be inherited as well. Nonetheless, our results broadly agree with [17], showing that when the effect of the background and effects of mutations within the QTL are large, adaptation proceeds without selective sweeps (Fig 4). We additionally show that the background explains considerable variation in many summary statistics, in particular those related to fixations and selective sweeps (Fig 3). Together with empirical observations of varying fitness effects for QTL in different backgrounds [46, 47, 48], our results highlight that evolutionary models of QTL cannot ignore the effects of genetic background.

### Genetic architecture of quantitative traits after adaptation

The genetic architecture of a trait is an important feature in the study of adaptation, influencing both the response to selection as well as the power to detect causal loci for a trait. Our two example traits show that different adaptation processes lead to different patterns of the genetic architecture matrix. Because Trait 1 only reached the new optimum shortly before we assess the the genetic architecture, the values are distributed asymmetrically around the zero effect size bin. Trait 2 reached the new optimum very early and therefore is more similar to an equilibrium genetic architecture with effect sizes close to zero at higher frequency and larger effect sizes at very low frequencies. These differences even between two highly correlated genetic architectures show that in addition to the input parameters, the time that passed since the new optimum was reached has an influence on the genetic architecture we observe in a population.

Using a machine-learning approach that trained on a subset of our simulations, we were able to identify the parameters that explained the largest proportion of variation among the genetic architectures studied (Fig 5). We found that demographic change plays a key role in determining the present genetic architecture, explaining as much as 55% (growth and bottleneck combined) of the variation we observed. For example, recent population growth leads to an increased number of low frequency mutations; this effect drives many of the observed differences between genetic architecture matrices of different demographies. We observed a high correlation (0.83–0.99) between genetic architectures with similar population demographies, suggesting that making inference about the process of adaptation from present day genetic architecture will have greater power in situations where the demography can be independently inferred. The result confirms the theoretical prediction that the combination of different allele frequency shifts at a large number of loci lead to similar trait architecture [49]. However when other statistics, such as information about fixations, effect size distributions observed in

present populations, number and type of selective sweeps and the demography are added as parameters to the modern genetic architecture, we can predict the evolutionary rate,  $\sigma_m$ , and  $V_S$  with 70% accuracy.

In addition to the effect of population growth, other input parameters do contribute substantially to variation in the genetic architecture, including the strength of stabilizing selection. Simons et al. [50] and [51] suggest that rare alleles are unlikely to contribute substantially to trait variance, but our models show that rare alleles can explain a large proportion of the variation when effect sizes are large. This is more consistent with the findings of [21], who showed that population growth leads to an increase proportion of genetic variance explained by rare alleles. The lack of consensus might result from differences in the models: while [50] models selection on fitness directly and [51] a quantitative trait under stabilizing selection with pleiotropy, our models and that of [21] consider selection on traits that are directly correlated to fitness.

### Maize domestication

Quantitative traits have been extensively studied in maize and breeders have made steady progress selecting traits for ever increasing trait values. But despite decades of observation that many important traits in maize are polygenic and work identifying QTL underlying domestication-related phenotypes [52], there has been little attention to the process of quantitative trait adaptation during maize domestication [but see 53, 23]. Nevertheless, many domestication traits, are polygenic and controlled by a number of loci with varying effect sizes [23]. Archaeological records of maize domestication traits show that adaptation took several thousand years [24]. Our example Trait 1 matches this pattern, representing an adaptation time of almost  $0.1N$  generations, equivalent to 10,000 years for a population similar to that of the wild ancestor of maize [22, Fig 6]. Trait 1 also leads to a reduction in genetic variance compared to the equilibrium population (wild ancestor), again matching observed data [23].

Trait 2, on the other hand, differs dramatically in a number of ways. It reached the new optimum extremely quickly, and diversity in the present is actually slightly higher than at the time of the optimum shift (Fig 6). The behavior of Trait 2 most closely resembles that of resistance traits with few large effect QTL [54]. We only look at the genetic variance of mutations that affect a single trait, but the overall diversity of a population is based on a combination of traits with different trait architectures and neutral parts of the genome. The observed reduction in genetic diversity of domesticates could partly be due to the distant optimum shift and partly, because of the population bottleneck experienced during domestication.

The difference in trait adaptation and genetic variance trajectory can be partially explained by the fixations and selective sweeps of beneficial alleles. The number of fixations revealed that as expected far more mutations fixed for Trait 1 than for Trait 2, as in Trait 1 many more sites are segregating in the equilibrium population, but the number of sweeps was also much higher. This is corrected for sites that fix due to genetic drift and shows that the larger relative distance to the new optimum changes the adaptation pattern. In maize it has been shown that domestication led to an accumulation of deleterious alleles, which so far was mainly attributed to the domestication bottleneck because no increase in deleterious alleles near major domestication genes was found [55]. For quantitative traits, small effect deleterious fixations could be distributed more uniformly across the genome and fix even without population bottlenecks. In general there are only few hard selective sweeps observed in maize and 84% of fitness related SNPs were already segregating in teosinte [56]. Our traits show that depending on the relative distance to the new optimum the type of selective sweeps changes. While sweeps come primarily from standing variation for traits that are close to the new optimum, for distant optima hard

sweeps are observed more frequently because standing variation is exhausted. The overall pattern of selective sweeps in the maize genome is a result of selection on a combination of traits and probably involves pleiotropic effects that can prevent fixation of new mutations even if they have large effects on a single trait [48].

### Signature of polygenic adaptation in genomic data

The recently suggested omnigenic model predicts that regulatory networks are sufficiently interconnected that many loci even outside the most biologically relevant genetic pathways can nonetheless affect a trait [57]. If indeed many traits are omnigenic, a quantitative evolutionary model as employed in our simulations is well suited for making inferences about observations in genomic data. Large sets of genomic and phenotypic data are becoming increasingly available, facilitating the study of the role of polygenic adaptation. Our results help to understand the implications of different parameters for the interpretation of such studies and provide targets for new selection tests that explicitly test for polygenic adaptation and the underlying genetic architecture. We show, for example, that selective sweeps can play a crucial role during polygenic adaptation and should be integrated into detection methods, as some approaches to investigate polygenic adaptation from shifts in allele frequencies may lose power if large effect alleles are fixed in the population in which effects are estimated [6, 39, 58].

Inferring polygenic adaptation and the underlying parameters in empirical data can provide important insight into the evolution of complex phenotypes. For experimental evolution scenarios in which the ancestral populations are known, the distance between the initial and the final optimum can be inferred from phenotype data, but for natural populations this may be more challenging. Our results indicate that the relative distance could be inferred from genomic data via estimates of the genetic architecture if the demographic history is known. One current challenge of transferring simulation results to empirical data is the computational limitation of simulating whole genome sequences in large populations. Faster implementations will allow simulation of larger regions and include neutral sites [59], and could be used to train machine learning models in order to predict the evolutionary history of a population from existing data coming from association studies. The implementation of machine learning trained on simulated data has been successfully applied to identify a number of population genetic patterns [60], and is a promising avenue for future work.

## Materials and methods

### Model

We simulated a quantitative trait under stabilizing selection with an optimum of 0 that adapted to a discrete optimum change to a value of 10. The population was diploid and mated randomly. Phenotypes followed a purely additive model in which the genotypic values at a given locus with an allele of effect size  $a$  were 0,  $0.5a$  and  $a$  for homozygous ancestral, heterozygous and homozygous derived genotypes. We modeled 20 QTL resembling 50kb regions, each with a 4 kb “genic” region centered in a 46 kb “intergenic” region. In the intergenic region mutations that affect the phenotype appeared with 1% probability of the genic region, leading to approximately 10% of mutations in intergenic regions and 90% in the 4kb genic regions. Starting with a neutral substitution rate of  $3 \times 10^{-8}$  per site per generation [61], we then assumed that only 10% of all mutations affect the trait of interest, resulting in a mutation rate of  $3 \times 10^{-9}$  per site per generation and a total per gamete mutation rate of  $3 \times 10^{-3}$  per generation. Regions were unlinked (50 cM distance), and within regions the recombination rate was  $5 \times 10^{-8}$  per site per generation (0.05 per gamete).

**Fitness.** We used a Gaussian fitness function in which an individual's fitness  $w$  was modeled as:

$$w = \exp\left[-\frac{(z - z_{opt})^2}{2V_s}\right] \quad (1)$$

where  $z$  is the trait value of an individual,  $z_{opt}$  is the population optimum trait value, and  $V_s$  modulates the possible deviation from the optimum. This standard model for traits under stabilizing selection is well suited for populations at equilibrium [26, chapter 7]. Under strong directional selection, however, this model greatly amplifies fitness differences among individuals in the tails of the phenotypic distribution. During the adaptive phase of the simulation, we calculated individual fitness following Eq 1, but then apply truncation selection by assigning a fitness of 1 to the top 50% of the distribution of  $w$  and 0 for the remaining 50%. This model allowed for truncation selection on  $z$ , while the population was distant from the new optimum, but allows for selection against phenotypes that surpass the new optimum during the final stages of adaptation. We stopped truncation selection once the population mean reached the new optimum, returning the population to stabilizing selection using fitness values calculated in Eq 1.

**Initial genetic variance.** The genetic variance at equilibrium can be approximated by the house of cards (HoC) model [12, 26]:

$$\mathbb{E}[V_G] = 4\mu V_s \quad (2)$$

and the stochastic HoC approximation [chapter 7, 26, 27]:

$$\mathbb{E}[V_G]_{SHC} = \frac{4\mu V_s}{1 + V_s/(N\sigma_m^2)} \quad (3)$$

We simulated five different values of  $V_s$  (1, 5, 10, 20, 50) to modulate the genetic variance of the equilibrium population.

**Effect size of new mutations.** We used a Gaussian distribution around zero for the effect size of new mutations and five different standard deviations ( $\sigma_m = 0.01, 0.05, 0.1, 0.3, 0.9$ ) to create traits with different effect sizes. Given a fixed optimum of 10, this distribution of effect sizes in combination with  $V_s$  effectively parameterize the distance to the new optimum, from a minimum distance of 11.5 z-scores (phenotypic standard deviations) to a maximum of 158.2 z-scores.

**Background.** Computational limitations do not allow simulation of an entire eukaryotic genome, so we added a heritable background ( $G_B$ ) to our simulations to account for the adaptive potential of the rest of the genome.

$$G_B \sim \mathcal{N}(G_{mp}, \sigma^2) \quad (4)$$

where  $G_B$  is the value of the genomic background of an individual,  $G_{mp}$  is the mid-parent genotypic value and  $\sigma^2$  is the variance of the parental trait values [62, chapter 9]. Hence,  $G_B$  is drawn from a normal distribution around the mid-parent value.

$$P = \psi \times G_B + (1 - \psi) \times G \quad (5)$$

The trait value of an individual  $P$  is then given by the sum of its genetic value  $G$  and the genomic background  $G_B$ , weighted by  $\psi$ , the proportion of trait variation represented by the background. We modeled four different background levels ( $\psi = 0, 0.1, 0.5, 0.95$ ).

**Demography.** To study the effect of population bottlenecks and expansion, we simulated a total of 12 different demographic scenarios with varying strength of a single bottleneck and

subsequent growth (S1 Fig). In scenarios with a bottleneck, an instantaneous reduction in population size occurs immediately after the burn-in and is followed by exponential growth over the length of the simulation ( $0.1 \times N_{anc}$  generations).

## Simulations

Using the above described parameters we simulated 100 replicates each of 25 different equilibrium traits using fwdpy11 v1.2a (<https://github.com/molpopgen/fwdpy11>), a Python package using the fwdpp library [63]. These 25 traits differed in their combination of  $V_S$  and  $\sigma_m$  and were run for a burn-in of  $10 N_{anc}$  generations (S3 Fig). Subsequently, each of the 1,200 parameter combinations was run for  $0.1 N_{anc}$  starting from these equilibrium traits.

We simulate a population of 10,000 individuals for 1,000 ( $0.1 N_{anc}$ ) generations after a burn-in of 100,000 generations to reach equilibrium.

The population mean trait values and variances were recorded every generation and entire populations, including individual trait values, mutations and effect sizes, were recorded every 10 generations for the first 100 generations after burn-in and then every 100 generations thereafter.

## Analysis

**Sweeps.** To identify selective sweeps, we used binomial sampling to simulate the sojourn time of neutral alleles arising in populations undergoing each of the demographic models described above. Mutations that were lost or that fixed before the end of the burn-in were ignored. We ran 10,000 replicates for each of the 12 demographic models and recorded the time it took a mutation that fixed within the last  $0.1N$  generations (similar to our selection model) to fix in this random model. These simulations provided a null distribution to which we compared selected mutations in our quantitative trait simulation (S2 Fig). We defined as a sweep any mutation that fixed faster than 99% of neutral alleles and categorized them as hard or from standing variation depending on whether the mutation arose before or after the optimum shift.

**Machine learning.** For each of the 120,000 simulations we calculated various summary statistics using the pandas version 0.21.0 and numpy version 1.12.1 Python libraries [64, 65]. These include statistics related to adaptation, selective sweeps, segregating sites, and fixed mutations; S1 Table. contains a full list of parameters used for prediction and importance estimation.

To identify the importance of input variables we trained a random forest and extracted the relative importance of the input parameters. We employed the RandomForestRegressor of sklearn 0.19.0 [66] with 100 trees to extract parameter importance by training the model using input parameters as features of the whole dataset and predicting a summary statistic. The prediction accuracy for all parameters was then estimated by 10-fold nested cross validation (training using 80% of the data) as well as root-mean-square deviation normalized by the range of values observed (NRMSD) implemented in sklearn [66], and the process repeated for each summary statistic of interest (S1 Table).

To compare the genetic basis of traits between scenarios we define the genetic architecture as the matrix of allele frequencies and effect sizes for each simulation. Allele frequencies were split into 7 discrete bins ( $0-10^{-4}$ ,  $10^{-4}-10^{-3}$ ,  $10^{-3}-10^{-2}$ ,  $10^{-2}-0.1$ ,  $0.1-0.5$ ,  $0.5-0.9$ ,  $0.9-1$ ) and effect sizes were split into 9 quantiles, as absolute effect sizes were strongly dependent on the input effect size. Relative occurrence frequencies (summing to 1 over the whole matrix) of segregating sites in each frequency-effect size combination were calculated for each simulation. These values were used to train a random forest model and extract parameter importance.

Parameter importance was estimated by predicting frequencies of each effect size bin from the input parameters. Prediction accuracy was again assessed by 10-fold cross validation implemented in sklearn [66]. Additionally, we calculated pairwise correlations of genetic architecture matrices in the final generation between all possible pairs of scenarios using the mean of all simulation replicates.

## Maize domestication

We took a closer look at two sets of simulations that represent diverging traits under a demographic model similar to that of maize domestication ( $N_{\text{bottleneck}} = 0.05 \times N_{\text{anc}}$ ;  $N_{\text{final}} = 10 \times N_{\text{anc}}$ ). For these simulations we assumed no genetic background ( $\psi = 0$ ). Trait 1 represents a trait with new mutations of small effect ( $\sigma_m = 0.01$ ) and strong stabilizing selection ( $V_S = 1$ ), while Trait 2 has new mutations of large effect ( $\sigma_m = 0.9$ ) and weaker stabilizing selection ( $V_S = 50$ ).

## Supporting information

**S1 Fig. Demographics.** Bottlenecks and growth models.

(TIF)

**S2 Fig. Detection of selective sweeps.** Distribution of fixation times from neutral single locus simulations (red) and forward simulations with selection (green). The grey area denotes the 99% confidence interval of neutral fixation time. Fixations outside the confidence interval are considered selective sweeps.

(TIF)

**S3 Fig. Genetic variance during burn-in.** The genetic variance in each generation over 10  $N_{\text{anc}}$  generations for each parameter set. Solid horizontal lines denote the House of Cards approximation of  $V_G$  [12]. Scenarios with small  $\sigma_m$  and large  $V_S$  do not reach the expected  $V_G$  because mutations are too small to “fill up” the variance volume. However, their equilibrium variance is well approximated by the stochastic House of Cards approximation [27, dashed line].

(TIF)

**S4 Fig. Equilibrium fitness.** Fitness for each burn-in parameter combination after 10N generations.

(TIF)

**S5 Fig. Relative parameter importance.** Relative parameter importance inferred by Random Forest machine learning for three parameter categories. 1) Adaptation, trait related parameters describing adaptation speed and potential for future adaptation. 2) Fixations, summary statistics for mutations that were fixed during trait adaptation and 3) segregating sites in the final generation of the simulations. Top panel indicating prediction accuracy as calculated by 10-fold nested cross validation and normalized relative mean squared error.

(TIF)

**S6 Fig. Fixations.** A) Total number of fixations B) Mean effect size of fixations C) Mean effect size of positive fixations D) Mean effect size of negative fixations.

(TIF)

**S7 Fig. Timing of selective sweeps for maize domestication simulations.** Shown is the relationship between effect size and the generation of fixation for mutations for Trait 1 (left,

$\sigma_m = 0.01$  and  $V_S = 1$ ) and Trait 2 (right,  $\sigma_m = 0.9$  and  $V_S = 50$ ).  
(TIF)

**S8 Fig. Fitness effects of mutations.** Fitness effects of mutations at the onset of directional selection (0.001–0.012N), before the new optimum is reached (0.001–0.012N) and after the new optimum has been reached (0.012–0.022N).

(TIF)

**S1 Table. Predicted summary statistics for feature importance estimation.**

(PDF)

## Acknowledgments

We would like to thank Patrick Thorwarth and Bernd Stetter for recommendations on machine learning, and CSI Davis, Graham Coop and Michael Turelli for helpful discussion. We are also grateful for the “wes anderson” color palettes that we employed for plots in this manuscript (<https://github.com/karthik/wesanderson>).

## Author Contributions

**Conceptualization:** Markus G. Stetter, Kevin Thornton, Jeffrey Ross-Ibarra.

**Data curation:** Markus G. Stetter.

**Formal analysis:** Markus G. Stetter.

**Funding acquisition:** Markus G. Stetter, Jeffrey Ross-Ibarra.

**Investigation:** Markus G. Stetter, Jeffrey Ross-Ibarra.

**Methodology:** Markus G. Stetter, Kevin Thornton, Jeffrey Ross-Ibarra.

**Project administration:** Markus G. Stetter.

**Resources:** Markus G. Stetter, Jeffrey Ross-Ibarra.

**Software:** Kevin Thornton.

**Supervision:** Jeffrey Ross-Ibarra.

**Validation:** Markus G. Stetter, Jeffrey Ross-Ibarra.

**Visualization:** Markus G. Stetter.

**Writing – original draft:** Markus G. Stetter, Jeffrey Ross-Ibarra.

**Writing – review & editing:** Markus G. Stetter, Kevin Thornton, Jeffrey Ross-Ibarra.

## References

1. Pritchard JK, Pickrell JK, Coop G. The genetics of human adaptation: hard sweeps, soft sweeps, and polygenic adaptation. *Current biology*. 2010; 20(4):R208–R215. <https://doi.org/10.1016/j.cub.2009.11.055> PMID: 20178769
2. Smith JM, Haigh J. The hitch-hiking effect of a favourable gene. *Genetics Research*. 1974; 23(1):23–35. <https://doi.org/10.1017/S0016672300014634>
3. Hermisson J, Pennings PS. Soft sweeps: molecular population genetics of adaptation from standing genetic variation. *Genetics*. 2005; 169(4):2335–2352. <https://doi.org/10.1534/genetics.104.036947> PMID: 15716498
4. Hermisson J, Pennings PS. Soft sweeps and beyond: understanding the patterns and probabilities of selection footprints under rapid adaptation. *Methods in Ecology and Evolution*. 2017; 8(6):700–716. <https://doi.org/10.1111/2041-210X.12808>

5. Pritchard JK, Di Rienzo A. Adaptation—not by sweeps alone. *Nature Reviews Genetics*. 2010; 11(10):665–667. <https://doi.org/10.1038/nrg2880> PMID: 20838407
6. Berg JJ, Zhang X, Coop G. Polygenic adaptation has impacted multiple anthropometric traits. *bioRxiv*. 2017; p. 167551.
7. Vignieri SN, Larson JG, Hoekstra HE. The selective advantage of crypsis in mice. *Evolution*. 2010; 64(7):2153–2158. <https://doi.org/10.1111/j.1558-5646.2010.00976.x> PMID: 20163447
8. Wallace J, Larsson S, Buckler E. Entering the second century of maize quantitative genetics. *Heredity*. 2014; 112(1):30–38. <https://doi.org/10.1038/hdy.2013.6> PMID: 23462502
9. Barton NH, Keightley PD. Understanding quantitative genetic variation. *Nature a-z index*. 2002; 3(1):11–21.
10. Wollstein A, Stephan W. Adaptive fixation in two-locus models of stabilizing selection and genetic drift. *Genetics*. 2014; 198(2):685–697. <https://doi.org/10.1534/genetics.114.168567> PMID: 25091496
11. Sanjak JS, Sidorenko J, Robinson MR, Thornton KR, Visscher PM. Evidence of directional and stabilizing selection in contemporary humans. *Proceedings of the National Academy of Sciences*. 2017; p. 201707227.
12. Turelli M. Heritable genetic variation via mutation-selection balance: Lerch’s zeta meets the abdominal bristle. *Theoretical population biology*. 1984; 25(2):138–193. [https://doi.org/10.1016/0040-5809\(84\)90017-0](https://doi.org/10.1016/0040-5809(84)90017-0) PMID: 6729751
13. Barton N. The maintenance of polygenic variation through a balance between mutation and stabilizing selection. *Genetics Research*. 1986; 47(3):209–216. <https://doi.org/10.1017/S0016672300023156>
14. Johnson T, Barton N. Theoretical models of selection and mutation on quantitative traits. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*. 2005; 360(1459):1411–1425. <https://doi.org/10.1098/rstb.2005.1667> PMID: 16048784
15. Jain K, Stephan W. Rapid adaptation of a polygenic trait after a sudden environmental shift. *Genetics*. 2017; 206(1):389–406. <https://doi.org/10.1534/genetics.116.196972> PMID: 28341654
16. Kopp M, Hermisson J. The genetic basis of phenotypic adaptation II: the distribution of adaptive substitutions in the moving optimum model. *Genetics*. 2009; 183(4):1453–1476. <https://doi.org/10.1534/genetics.109.106195> PMID: 19805820
17. Chevin LM, Hospital F. Selective Sweep at a Quantitative Trait Locus in the Presence of Background Genetic Variation. *Genetics*. 2008; 180(3):1645–1660. <https://doi.org/10.1534/genetics.108.093351> PMID: 18832353
18. Lande R. The response to selection on major and minor mutations affecting a metrical trait. *Heredity*. 1983; 50(1):47–65. <https://doi.org/10.1038/hdy.1983.6>
19. de Vladar HP, Barton N. Stability and response of polygenic traits to stabilizing selection and mutation. *Genetics*. 2014; 197(2):749–767. <https://doi.org/10.1534/genetics.113.159111> PMID: 24709633
20. Pavlidis P, Metzler D, Stephan W. Selective sweeps in multilocus models of quantitative traits. *Genetics*. 2012; 192(1):225–239. <https://doi.org/10.1534/genetics.112.142547> PMID: 22714406
21. Lohmueller KE. The impact of population demography and selection on the genetic architecture of complex traits. *PLoS genetics*. 2014; 10(5):e1004379. <https://doi.org/10.1371/journal.pgen.1004379> PMID: 24875776
22. Beissinger TM, Wang L, Crosby K, Durvasula A, Hufford MB, Ross-Ibarra J. Recent demography drives changes in linked selection across the maize genome. *Nature Plants*. 2016; 2:16084. <https://doi.org/10.1038/nplants.2016.84> PMID: 27294617
23. Xue S, Bradbury PJ, Casstevens T, Holland JB. Genetic architecture of domestication-related traits in maize. *Genetics*. 2016; 204(1):99–113. <https://doi.org/10.1534/genetics.116.191106> PMID: 27412713
24. Benz BF, Cheng L, Leavitt SW, Eastoe C. El Riego and early maize agricultural evolution. *Histories of maize*. 2006; p. 73–82.
25. Bürger R. Linkage and the maintenance of heritable variation by mutation-selection balance. *Genetics*. 1989; 121(1):175–184. PMID: 2917713
26. Bürger R. The mathematical theory of selection, recombination, and mutation. vol. 228. Wiley Chichester; 2000.
27. Bürger R. Mutation-selection balance and continuum-of-alleles models. *Mathematical biosciences*. 1988; 91(1):67–83. [https://doi.org/10.1016/0025-5564\(88\)90024-7](https://doi.org/10.1016/0025-5564(88)90024-7)
28. Moran EV, Alexander JM. Evolutionary responses to global change: lessons from invasive species. *Ecology Letters*. 2014; 17(5):637–649. <https://doi.org/10.1111/ele.12262> PMID: 24612028
29. Turner TL, Stewart AD, Fields AT, Rice WR, Tarone AM. Population-based resequencing of experimentally evolved populations reveals the genetic basis of body size variation in *Drosophila melanogaster*. *PLoS genetics*. 2011; 7(3):e1001336. <https://doi.org/10.1371/journal.pgen.1001336> PMID: 21437274

30. Dudley J. From means to QTL: The Illinois long-term selection experiment as a case study in quantitative genetics. *Crop Science*. 2007; 47(Supplement\_3):S–20. <https://doi.org/10.2135/cropsci2007.04.0003IPBS>
31. Crow JF, Kimura M. Efficiency of truncation selection. *Proceedings of the National Academy of Sciences*. 1979; 76(1):396–399. <https://doi.org/10.1073/pnas.76.1.396>
32. Oz T, Guvenek A, Yildiz S, Karaboga E, Tamer YT, Mumcuyan N, et al. Strength of selection pressure is an important parameter contributing to the complexity of antibiotic resistance evolution. *Molecular biology and evolution*. 2014; 31(9):2387–2401. <https://doi.org/10.1093/molbev/msu191> PMID: 24962091
33. Hoekstra HE, Hoekstra JM, Berrigan D, Vignieri SN, Hoang A, Hill CE, et al. Strength and tempo of directional selection in the wild. *Proceedings of the National Academy of Sciences*. 2001; 98(16):9157–9160. <https://doi.org/10.1073/pnas.161281098>
34. Carlborg Ö, Haley CS. Epistasis: too often neglected in complex trait studies? *Nature Reviews Genetics*. 2004; 5(8):618. <https://doi.org/10.1038/nrg1407> PMID: 15266344
35. Jiang Y, Schmidt RH, Zhao Y, Reif JC. A quantitative genetic framework highlights the role of epistatic effects for grain-yield heterosis in bread wheat. *Nature genetics*. 2017; 49(12):1741. <https://doi.org/10.1038/ng.3974> PMID: 29038596
36. Polderman TJ, Benyamin B, De Leeuw CA, Sullivan PF, Van Bochoven A, Visscher PM, et al. Meta-analysis of the heritability of human traits based on fifty years of twin studies. *Nature genetics*. 2015; 47(7):702. <https://doi.org/10.1038/ng.3285> PMID: 25985137
37. Zhu Z, Bakshi A, Vinkhuyzen AA, Hemani G, Lee SH, Nolte IM, et al. Dominance genetic variation contributes little to the missing heritability for human complex traits. *The American Journal of Human Genetics*. 2015; 96(3):377–385. <https://doi.org/10.1016/j.ajhg.2015.01.001> PMID: 25683123
38. Mackay TF. Epistasis and quantitative traits: using model organisms to study gene–gene interactions. *Nature Reviews Genetics*. 2014; 15(1):22. <https://doi.org/10.1038/nrg3627> PMID: 24296533
39. Forsberg SK, Bloom JS, Sadhu MJ, Sadhu MJ, Kruglyak L, Carlborg Ö. Accounting for genetic interactions improves modeling of individual quantitative trait phenotypes in yeast. *Nature genetics*. 2017; 49(4):497. <https://doi.org/10.1038/ng.3800> PMID: 28250458
40. Charlesworth D, Willis JH. The genetics of inbreeding depression. *Nature reviews genetics*. 2009; 10(11):783. <https://doi.org/10.1038/nrg2664> PMID: 19834483
41. Yang J, Mezouk S, Baumgarten A, Buckler ES, Guill KE, McMullen MD, et al. Incomplete dominance of deleterious alleles contributes substantially to trait variation and heterosis in maize. *PLoS genetics*. 2017; 13(9):e1007019. <https://doi.org/10.1371/journal.pgen.1007019> PMID: 28953891
42. Burger R, Lynch M. Evolution and extinction in a changing environment: a quantitative-genetic analysis. *Evolution*. 1995; p. 151–163. <https://doi.org/10.1111/j.1558-5646.1995.tb05967.x> PMID: 28593664
43. Haldane JBS. A mathematical theory of natural and artificial selection, part V: selection and mutation. In: *Mathematical Proceedings of the Cambridge Philosophical Society*. vol. 23. Cambridge University Press; 1927. p. 838–844.
44. Wright S. Evolution in Mendelian populations. *Genetics*. 1931; 16(2):97–159. PMID: 17246615
45. Berg JJ, Coop G. A coalescent model for a sweep of a unique standing variant. *Genetics*. 2015; 201(2):707–725. <https://doi.org/10.1534/genetics.115.178962> PMID: 26311475
46. Symonds VV, Godoy AV, Alconada T, Botto JF, Juenger TE, Casal JJ, et al. Mapping quantitative trait loci in multiple populations of *Arabidopsis thaliana* identifies natural allelic variation for trichome density. *Genetics*. 2005; 169(3):1649–1658. <https://doi.org/10.1534/genetics.104.031948> PMID: 15654092
47. Doebley J, Stec A, Gustus C. *teosinte branched1* and the origin of maize: evidence for epistasis and the evolution of dominance. *Genetics*. 1995; 141(1):333–346. PMID: 8536981
48. Stitzer MC, Ross-Ibarra J. Maize domestication and gene interaction. *PeerJ Preprints*. 2018; 6:e26502v1.
49. Lynch M, Walsh B. *Genetics and analysis of quantitative traits*. vol. 1. Sinauer Sunderland, MA; 1998.
50. Simons YB, Turchin MC, Pritchard JK, Sella G. The deleterious mutation load is insensitive to recent population history. *Nature genetics*. 2014; 46(3):220. <https://doi.org/10.1038/ng.2896> PMID: 24509481
51. Simons YB, Bullaughey K, Hudson RR, Sella G. A population genetic interpretation of GWAS findings for human quantitative traits. *PLoS biology*. 2018; 16(3):e2002985. <https://doi.org/10.1371/journal.pbio.2002985> PMID: 29547617
52. Briggs WH, McMullen M, Gaut B, Doebley J. Linkage mapping of domestication loci in a large maize-teosinte backcross resource. *Genetics*. 2007; . <https://doi.org/10.1534/genetics.107.076497>

53. Brown PJ, Upadaya N, Mahone GS, Tian F, Bradbury PJ, Myles S, et al. Distinct genetic architectures for male and female inflorescence traits of maize. *PLoS genetics*. 2011; 7(11):e1002383. <https://doi.org/10.1371/journal.pgen.1002383> PMID: 22125498
54. Poland JA, Bradbury PJ, Buckler ES, Nelson RJ. Genome-wide nested association mapping of quantitative resistance to northern leaf blight in maize. *Proceedings of the National Academy of Sciences*. 2011; 108(17):6893–6898. <https://doi.org/10.1073/pnas.1010894108>
55. Wang L, Beissinger TM, Lorant A, Ross-Ibarra C, Ross-Ibarra J, Hufford MB. The interplay of demography and selection during maize domestication and expansion. *Genome biology*. 2017; 18(1):215. <https://doi.org/10.1186/s13059-017-1346-4> PMID: 29132403
56. Swarts K, Gutaker RM, Benz B, Blake M, Bukowski R, Holland J, et al. Genomic estimation of complex traits reveals ancient maize adaptation to temperate North America. *Science*. 2017; 357(6350):512–515. <https://doi.org/10.1126/science.aam9425> PMID: 28774930
57. Boyle EA, Li YI, Pritchard JK. An expanded view of complex traits: from polygenic to omnigenic. *Cell*. 2017; 169(7):1177–1186. <https://doi.org/10.1016/j.cell.2017.05.038> PMID: 28622505
58. Crawford NG, Kelly DE, Hansen ME, Beltrame MH, Fan S, Bowman SL, et al. Loci associated with skin pigmentation identified in African populations. *Science*. 2017; 358(6365):eaan8433. <https://doi.org/10.1126/science.aan8433> PMID: 29025994
59. Kelleher J, Thornton K, Ashander J, Ralph P. Efficient pedigree recording for fast population genetics simulation. *bioRxiv*. 2018; p. 248500.
60. Schrider DR, Kern AD. Supervised Machine Learning for Population Genetics: A New Paradigm. *Trends in Genetics*. 2018; <https://doi.org/10.1016/j.tig.2017.12.005> PMID: 29331490
61. Clark RM, Tavaré S, Doebley J. Estimating a nucleotide substitution rate for maize from polymorphism at a major domestication locus. *Molecular biology and evolution*. 2005; 22(11):2304–2312. <https://doi.org/10.1093/molbev/msi228> PMID: 16079248
62. Falconer DS. *Introduction to quantitative genetics*. Oliver And Boyd; Edinburgh; London; 1960.
63. Thornton KR. A C++ template library for efficient forward-time population genetic simulation of large populations. *Genetics*. 2014; 198(1):157–166. <https://doi.org/10.1534/genetics.114.165019> PMID: 24950894
64. McKinney W. Data Structures for Statistical Computing in Python. In: van der Walt S, Millman J, editors. *Proceedings of the 9th Python in Science Conference*; 2010. p. 51–56.
65. Walt Svd, Colbert SC, Varoquaux G. The NumPy array: a structure for efficient numerical computation. *Computing in Science & Engineering*. 2011; 13(2):22–30. <https://doi.org/10.1109/MCSE.2011.37>
66. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*. 2011; 12:2825–2830.