

# Demography and linked selection in wild and domesticated maize

Timothy M. Beissinger<sup>\* † ‡</sup>, Li Wang<sup>§</sup>, Kate Crosby<sup>\*</sup>, Arun Durvasula<sup>\*</sup>, Matthew Hufford<sup>§</sup>, and Jeffrey Ross-Ibarra<sup>\* ¶</sup>

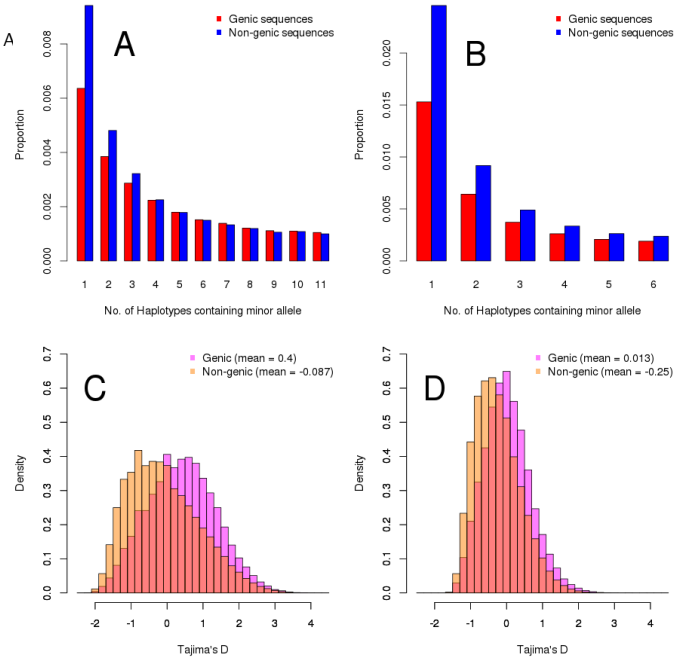
<sup>\*</sup>Dept. of Plant Sciences, University of California, Davis, CA, USA, <sup>†</sup>US Department of Agriculture, Agricultural Research Service, Columbia, MO, USA, <sup>‡</sup>Division of Plant Sciences, University of Missouri, Columbia, MO, USA, <sup>§</sup>Iowa State University, Ames, IA, USA, and <sup>¶</sup>Genome Center and Center for population biology, University of California, Davis, CA, USA

Submitted to Proceedings of the National Academy of Sciences of the United States of A

Unique selective and demographic operate on domesticated plant species. These forces interact during and after domestication to generate the patterns of DNA variability that are persistent in crop species today. To quantify the interplay between demography and selection, we investigated genetic diversity in maize, one of the most important crops for food, feed, and fuel world-wide. We utilized whole genome sequence data from 23 maize and 13 teosinte individuals to make inferences. We obtained a complete estimate of the population size fluctuations and other demographic parameters experienced by maize as it was domesticated from teosinte. Here, we show that maize went through a domestication bottleneck with a population size of approximately 5% that of teosinte before it experienced rapid population size expansion post-domestication. We observe that hard sweeps, specifically positive selection on new genic mutations, are not the primary force driving maize evolution. We find that a reduced population size during domestication decreased the efficiency of purifying selection to purge deleterious alleles from maize. However, expansion after domestication has since increased the efficiency of purifying selection to levels superior to those seen in teosinte. Our results demonstrate that in domesticated species, demographic and selective history in the ancient and recent past both contribute to genetic variability that is present today, providing substantial implications for the continued improvement of domesticated species.

Domesticated plant species evolve in a unique fashion compared to their wild counterparts [1]. This is a result of both the anthropomorphic nature of artificial selection on domesticates [2] as well as the demographic characteristics of the domestication bottleneck(s) that they tend to have experienced [3]. However, the complex interplay between selective pressures and demographic limitations, and the impact that this interplay has on identifying selection and understanding demography, is not fully understood. Although a large body of research that involves searching the genomes of domesticated species for evidence of positive selection exists [4–7], these studies tend to focus on identifying or mapping particular genes or regions that play an important role in phenotypic evolution. In contrast, knowledge regarding the impacts that demography and selection have on whole-genome patterns of genetic variability remains limited.

For instance, supposed neutrally-evolving DNA is often used to estimate historical demographic parameters of a population such as effective population size, structure, and expansion history [8,9]. However, researchers have called into question whether or not there are reasonable approaches for identifying neutral regions of the genome, since the effects of selection can be wide-ranging [10,11]. For example, in *Drosophila* it appears that the majority of the genome is impacted by the effects of selection through linkage [12]. A natural next question, therefore, is how can demographic parameters be estimated independently of selective parameters? Human researchers have attempted to address this problem by limiting analyses to only sites far from genes [13], but as *Drosophila* demonstrates, it is difficult to be certain that sites even in gene-poor regions of the genome are not influenced by linked selection. Ultimately, an understanding of which sites have



**Fig. 1.** The folded SFS and Tajimas D were calculated both inside and outside of genes for maize and teosinte. In both taxa, a dearth of rare alleles was observed within genes relative to outside of genes, which led to higher values of Tajima’s D outside of genes than inside. **A:** Folded SFS for maize; **B:** Folded SFS for teosinte; **C:** Maize Tajima’s D; **D:** Teosinte Tajima’s D. skip the SFS and make this a panel plot of pi and TajD in maize and teosinte

## Significance

Patterns of linked selection, or the impact of selection on sites that neighbor a functional variant, have been carefully evaluated in only a few species. not really. we need to rephrase here and cite the Corbett-Detig plos-bio paper somewhere – they looked at a lot of species including maize! In this work, we demonstrate that selection against deleterious mutations leaves a pronounced signature on the maize genome, reducing diversity in and immediately around genes. We show how demography interacts with selection to impact genome-wide patterns of diversity, including the important observation that rapid population expansion can increase the efficiency of selection as much as a sudden population bottleneck can weaken it. Along the way, we develop the first estimate the demographic parameters of the maize domestication from whole genome sequence data.

## Reserved for Publication Footnotes

wide patterns of variability through linkage, is required for reasonable demographic inference.

Maize represents an excellent organism to study these phenomena. Maize is a species of tremendous importance worldwide as both a staple crop [14] and as a model for understanding plant evolution [15]. Broadly speaking, archaeological and genetic studies have established that maize domestication is likely to have taken place in Mexico approximately 9,000 years bp [16,17]. Teosinte, the most recent wild ancestor to maize, remains extant throughout much of the Americas [18]. Additionally, several large-effect domestication loci [19–21] and putative domestication regions [4] have been identified. But despite all that is known about maize domestication, the parameters of the domestication process remain uncertain. Specifically, the size of the maize domestication bottleneck has not been estimated independently of the bottleneck’s duration, nor are there sequence-based estimates of the effective population size of modern maize. Sequence information from maize and teosinte plants may therefore be utilized to address these questions.

To that end, the objectives of our study were to 1) investigate the relative importance of different forms of selection on whole-genome variability in both maize and teosinte 2) research the impact that the domestication process has had on genetic variability in maize, and how this compares to the impact of a different demographic history in teosinte; and 3) precisely estimate the parameters of the maize domestication bottleneck. We show that our third objective, estimating the parameters of domestication, is not possible without first completing objectives one and two, which demonstrate that as in humans [13], the majority of maize non-genic DNA may reasonably be treated as neutral. We achieve these objectives by utilizing whole-genome-sequence information from 23 maize and 13 teosinte lines sequenced as part of the Maize HapMap 2 panel [22].

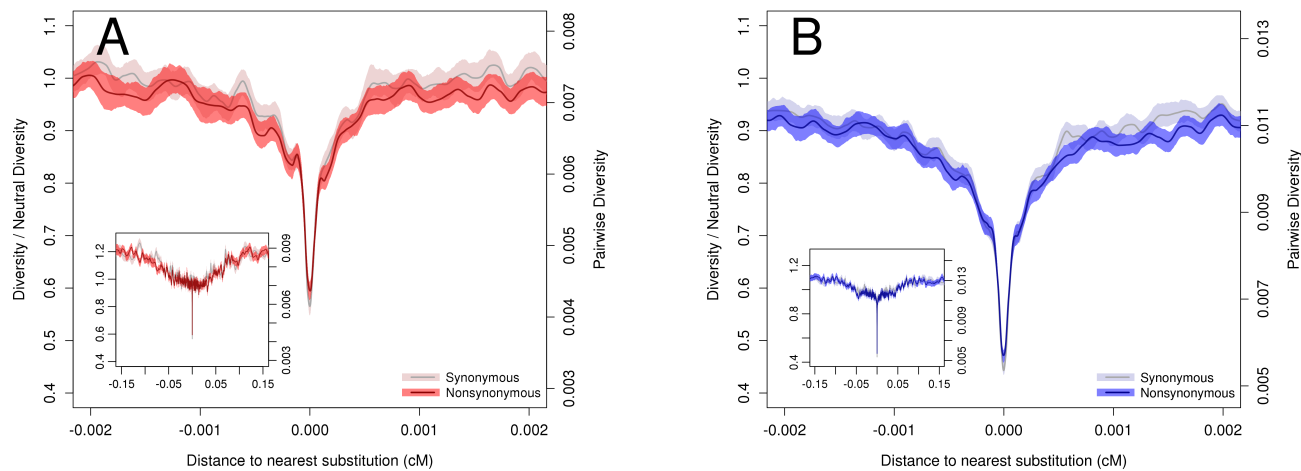
## Results

**Patterns of diversity differ between genic and non-genic regions of the genome.** Our reanalysis of the maize HapMap 2 data [22] is consistent with earlier findings that patterns

of diversity differ between genic and non-genic regions of the genome in both maize and teosinte [4] (Figure 1). In maize, mean pairwise diversity ( $\pi$ ) within genes was significantly lower than at positions at least 5 kb away from genes (0.00668 vs 0.00691,  $p < 2 \times 10^{-44}$ ). Diversity differences in teosinte are even more pronounced (0.0088 vs. 0.0115,  $p \approx 0$ ). Differences were also apparent in the site frequency spectrum, with mean Tajima’s D positive in genic regions in both maize (0.4) and teosinte (0.013) but negative outside of genes (-0.087 in maize and -0.25 in teosinte,  $p < X$  for both comparisons). These observations suggest that diversity in genes is not evolving neutrally, but instead reduced by the impacts of selection on linked sites.

**Hard sweeps do not explain diversity differences.** Selection acting to increase the frequency of a new beneficial mutation will leave a signature of reduced diversity at surrounding linked sites **CITE**. To evaluate whether patterns of such “hard sweeps” could explain observed differences in diversity between genic and non-genic regions of the genome, we compared diversity around missense and synonymous substitutions between *Tripsacum* and either maize or teosinte (Figure 2). If a proportion of missense mutations have been fixed due to hard sweeps, diversity around these substitutions should be lower than around synonymous substitutions. We observe this pattern around the causative amino acid substitution in the domestication locus *tga1* (Figure S1), likely the result of a hard sweep during domestication [?, 21].

Genome-wide, however, we observe no differences in diversity between synonymous and missense substitutions in either maize or teosinte. Previous analyses in humans have suggested that this approach is limited by fact that a higher proportion of nonsynonymous substitutions will be found in genes under weak purifying selection and thus higher genetic diversity [27]. To address this concern we took advantage of a measure of evolutionary constraint using genomic evolutionary rate profile (GERP) scores [28] calculated across the maize genome [30]. We re-analyzed substitutions in subsets of genes with the highest and lowest 10% quantile of mean GERP score, putatively representing genes under the strongest and weakest purifying selection (Figure S2). As expected, we see higher overall



**Fig. 2.** Pairwise diversity surrounding synonymous and non-synonymous substitutions in **A** maize and **B** teosinte. Axes show both absolute diversity values and values relative to mean nucleotide diversity in windows  $\geq 0.01cM$  from a substitution. **is this correct or did i confuse things?** Lines depict a loess curve (span of 0.01) and shading represents bootstrap-based 95% confidence intervals. Inset plots depict a larger range on the x-axis.

diversity around substitutions in genes under weak purifying selection, but we still see no difference between synonymous and missense substitutions in either subset of the data. Taken together, these data suggest hard sweeps do not play a major role in patterning genic diversity in either maize or teosinte.

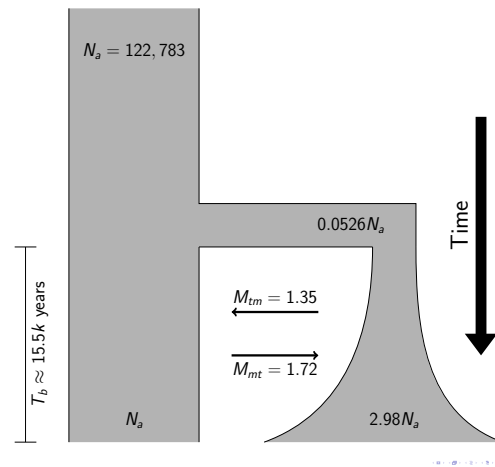
#### Pairwise diversity is heavily influenced by purifying selection.

Purifying selection refers to the situation where deleterious mutations arising in a population are continuously selected against. When this form of selection is operating, it can serve to reduce genetic variability at linked neutral sites, a phenomenon often called background selection [31]. Purifying and background selection lead to lower diversity within genes and other functional sites relative to neutral regions. We investigated purifying selection in maize and teosinte by evaluating the average magnitude of reduced diversity within genes and recovery away from genes in both taxa. To enable comparisons between maize and teosinte, we standardized by neutral levels within each taxa, where neutral levels were defined as mean diversity at positions at least 0.01 cM distal from genes. For both maize and teosinte, pairwise diversity was reduced within genes and a recovery of diversity away from genes was observed. Interestingly, a stronger reduction of diversity and slower recovery was observed for teosinte than for maize, implying that purifying selection has left a more pronounced signature on pairwise diversity in the teosinte genome (Figure 4).

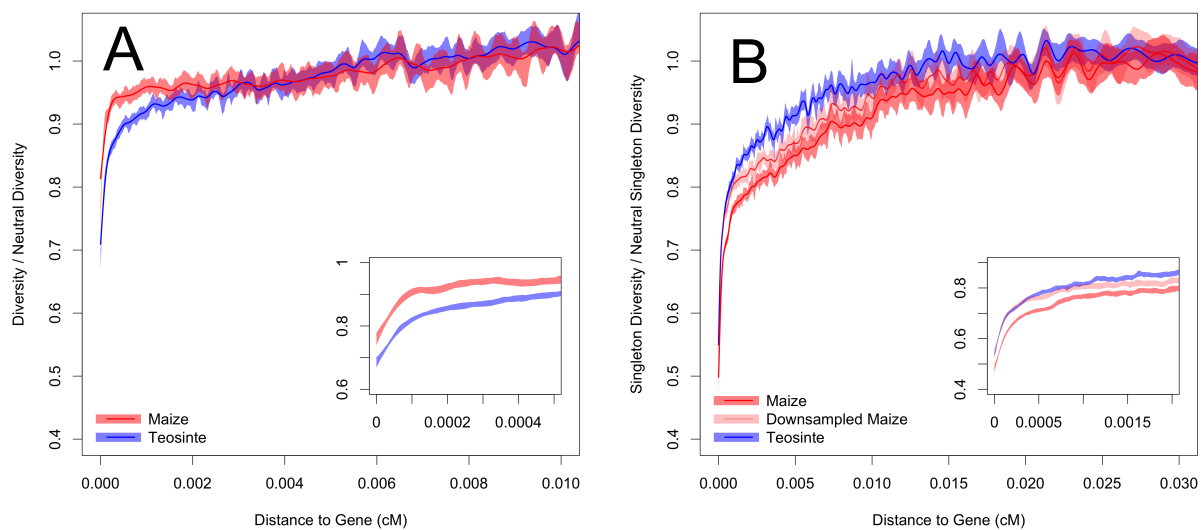
**Demography of maize domestication.** To explore whether differences in purifying selection between maize and teosinte can be explained by demographic processes, we first estimated the parameters of maize domestication using large-scale sequencing data involving 23 inbred maize landraces and 13 teosinte inbred lines included in the HapMap 2 panel [22]. The maize lines were collected from across the Americas and the teosinte lines came from central Mexico. Before estimating demography, we compared the site frequency spectrum (SFS) in genic and non-genic regions, and observed substantial differences

in the evolution of these classes of sites. For both maize and teosinte, the SFS within genes showed a dearth of low-frequency alleles and Tajima’s  $D$  [32] was therefore shifted to more positive values, as previously shown (Figure 4). This is consistent with the aforementioned purifying selection and indicates that genic regions are not evolving neutrally. Therefore, for demographic modeling we restricted analysis to sites at least 5 kb from gene start/stop positions.

We used diffusion approximation as implemented in  $\delta\alpha\delta i$  [9] to find the domestication parameters that best explain the joint site frequency spectrum of maize and teosinte. The model we optimized began with an equilibrium ancestral population of size  $N_a$  splitting into separate maize and teosinte



**Fig. 3.** Parameters of domestication as estimated by  $\delta\alpha\delta i$ . Notably, the maize effective population size ( $N_e$ ) during the domestication bottleneck appears to have consisted of approximately 5.26% the ancestral  $N_e$ , before recovering two at least 2.98 times as large as the ancestral  $N_e$ .  $M_{mt}$  and  $t_m$  are the reverse of what is described in text.



**Fig. 4.** Relative diversity versus distance to nearest gene in maize and teosinte. Shown are **A** pairwise nucleotide diversity and **B** singleton diversity. Relative diversity is calculated compared to the mean diversity in windows  $\geq 0.01cM$  or  $\geq 0.02cM$  from the nearest gene for pairwise diversity and singletons, respectively. Lines depict a loess curve (span of 0.01) and shading represents bootstrap-based 95% confidence intervals. Inset plots depict a smaller range on the x-axis.

populations  $T_B$  generations in the past. Moving forward in time from the split, teosinte maintained the ancestral population size, while maize experienced an immediate effective population size change to  $N_b$  individuals, followed by exponential growth to size  $N_m$ . Additionally, since the split  $M_{tm}$  individuals migrated from teosinte into maize and  $M_{mt}$  individuals migrated from maize to teosinte.

The most likely model suggested an ancestral  $\theta$  ( $4N\mu$ ) per base pair of 0.014734, which approximately matches a previous independent estimate [25]. Assuming a mutation rate of  $\mu = 3 \times 10^{-8}$  [33], this suggests an effective population size of  $N_a = 122,783$  individuals, a population split  $T_B = 15,523$  generations in the past, a maize bottleneck effective population size of  $N_b = 6,455$  individuals (5.26% of  $N_a$ ), a modern maize effective population size of  $N_m = 366,973$  individuals,  $M_{tm} = 1.35$  migrants per generation from teosinte to maize, and  $M_{mt} = 1.72$  migrants per generation from maize to teosinte (Figure 3). We note that a population split over 15 thousand generations before present precedes estimates from archaeological data, which suggest maize domestication began no more than 10,000 years before present [16]. This could reflect teosinte population structure that was present before domestication. Also, note that the genetic time of the population split must precede physiological changes that could be identified archaeologically. Additionally, because recent expansion is most evidenced by rare alleles, and since our limited sample size provides low power to detect rare alleles, we expect that the estimate of  $N_m = 366,973$ , or  $\sim 3N_a$ , is likely an underestimate.

Therefore, we utilized a complementary data set of 4,021 maize land-race individuals collected from across the Americas (SeeDs reference here) to obtain a less downward-biased estimate of the current maize effective population size. According to a singleton-based estimate of  $\theta$  [34], which was utilized since rapid post-bottleneck expansion implies that rare alleles will most accurately reflect current maize demography, we obtained the estimate  $N_e = 992,713.5$ . Although less biased than the previously mentioned estimate, we note that ascertainment bias of the genotyping platform used to generate these data again biases this estimate downward.

**Post-domestication expansion is apparent in patterns of singleton diversity.** Motivated by the rapid post-domestication expansion in maize that is evident from our demographic analysis, we conducted an analysis of genome-wide patterns of singleton diversity. Singleton diversity measures the abundance of alleles that are only observed in one individual among the sample. As a class, singleton alleles depict the most recent patterns of evolution, but they also have the lowest effect on pairwise diversity. Therefore, unlike pairwise diversity, patterns of singleton diversity reflect recent patterns of evolution. Since our sample size for maize ( $n=23$ ) was larger than teosinte ( $n=13$ ), maize singletons were analyzed directly as well as down-sampled to reflect the sample size of teosinte. For singletons, our definition of neutral diversity was defined as diversity at positions at least 0.02 cM distal from genes. When evaluating the data in this manner, a very different pattern emerged compared to what we saw with respect to pairwise diversity. Maize singleton diversity was at least as reduced as teosinte singleton diversity near genes, but recovered more slowly (Figure 4), implying that in the recent past maize has been more influenced by purifying selection than has teosinte.

Together, these findings suggest that demographic history has a strong influence on the effect of purifying selection. Historically, teosinte has had a larger population size than maize,

and only recently has maize population size overcome that of teosinte. Since the efficacy of purifying selection scales with population size, these results likely reflect changes in  $N_e$  more than they reflect underlying changes in selection pressure.

Similarly, to complement our investigation of pairwise diversity surrounding substitutions, we studied patterns of singleton diversity surrounding synonymous and non-synonymous substitutions in maize (Figure S3). A nearly identical pattern to that shown in Figure 2 was observed. This further demonstrates that even the most recent selection patterns in maize are not dominated by hard sweeps.

**Purifying selection results are robust to soft-sweeps.** As a category, instances of soft-sweeps, or cases of selection on alleles that have already reached intermediate frequency before becoming advantageous, are difficult to identify [35]. However, this form of selection is likely to be frequent during crop domestication [36]. We were concerned that the disparate patterns of pairwise and singleton diversity in maize and teosinte, which seem to reflect differences in the efficiency of purifying selection due to demographic history, may instead reflect differences in the location and abundance of soft-sweeps. Therefore, we re-analyzed patterns of pairwise and singleton diversity in maize and teosinte including only a subset of sites that are least likely to correspond to soft-sweeps. To achieve this, we performed a genome-wide scan for soft or hard sweeps in maize or teosinte based on the H12 statistic [37]. H12 is sensitive to both hard and soft sweeps. Sites nearest to genes that displayed among the top 20% of H12 values in either maize or teosinte were removed from our data set, and both pairwise and singleton diversity at remaining sites, as a function of the distance to the nearest gene, was calculated (Figure S4). This subset of the data is no longer subject to the potentially confounding effects of soft-sweeps, and nearly identical patterns of pairwise and singleton diversity surrounding genes are observed. This provides further evidence that the patterns of diversity we observe in maize as compared to teosinte are the result of rapid population expansion as opposed to differing selective patterns in the domesticated and wild taxa.

## Discussion

**Little positive selection on new genic mutations.** mention divergence time from *Tripsacum*, roughly 1-1.2 million years before present [26]. Our findings indicate that hard selective sweeps have not contributed substantially to genome-wide patterns of diversity in maize. More specifically, we observe that positive selection on new genic mutations in maize is not common compared to other drivers of diversity. This finding does not exclude the possibility that hard sweeps have taken place at non-genic sites, such as in those in enhancer regions. One known example of positive selection on a non-genic mutation involves the *tb1* locus of maize, which is one of the best characterized examples of positive selection on a “domestication gene” in any crop [38]. However, the maize *tb1* allele was already present in teosinte before domestication [39], so this this example is in agreement with our finding that hard sweeps are rare. The *gt1* locus is another well-characterized case of positive selection operating on standing variation at an enhancer region [20]. Instances of selection on standing variation such as these, often called soft sweeps, may be a major contributor to maize patterns of diversity [40], and this could be part of the explanation as to why hard sweeps appear to be so rare, despite obvious morphological differences between maize and teosinte.

Unfortunately, our ability to accurately identify soft sweeps, and particularly to distinguish them from hard sweeps, re-



mains limited [36,41]. However, we implemented a scan based on the H12 statistic, which is designed to identify both hard and soft sweeps with reasonable power [37]. Although the goal of this study was not to identify specific sweeps, it may be that the outlier sites distinguished by H12 are primarily composed of soft sweeps instead of hard. Similarly, previous studies scanning for evidence of positive selection in maize [4] may have picked up primarily evidence of soft sweeps. In domesticated species, artificial selection beginning at the onset of domestication represents a drastic shift in selective pressure. One can easily imagine a scenario in which previously neutral or slightly deleterious variants that were segregating in the population when this shift took place suddenly became beneficial, leading to soft sweeps. As appealing as this explanation is, however, our data do not speak to the abundance of soft sweeps, beyond demonstrating that they are one possibility for the lack of observed hard sweeps.

We should additionally note that our observations do not exclude the possibility of infrequent hard sweeps having taken place during maize evolution. In fact, the maize locus *tga1* that has been shown to correspond to a hard sweep at an amino-acid changing mutation [21]. Surrounding *tga1*, our data demonstrate a pattern consistent with a hard sweep. But, our data demonstrate that instances such as this are infrequent. This contrasts sharply with *Drosophila* [42] and *Capsella* [43], where differences in diversity surrounding synonymous and non-synonymous substitutions are clear suggesting hard-sweeps are abundant. However, it agrees with what has been seen for humans [44].

**Demography of domestication.** add note that demography inference will require regions far from genes, compare to [23–25]. cite some theory or empirical paper showing this to be the case in humans/drosophila Since we observed that classical hard sweeps do not explain patterns of maize diversity, it is perhaps unsurprising that we instead found evidence of demographic history contributing substantially to genetic diversity. Our estimate of the demographic history of maize, which is the first estimate of maize demographic history to utilize whole-genome-sequence (WGS) data, allowed us to disentangle the confounding parameters of bottleneck strength and bottleneck intensity. This confounding limited previous domestication estimates [23]. By utilizing the two-dimensional site frequency spectrum (SFS) between maize and teosinte for inference [9] we estimate that the maize effective population size during the bottleneck was approximately 5% that of teosinte. A bottleneck of this magnitude may have enabled moderately deleterious alleles to reach intermediate frequency in maize, while also increasing the probability that strongly deleterious alleles segregating at low frequency in teosinte were purged from the maize gene pool, due to elevated drift as is predicted during a domestication bottleneck [32].

After the initial bottleneck, our analysis suggests that maize  $N_e$  expanded to reach levels much greater than that of teosinte. Our WGS-based estimate shows that the  $N_e$  of landrace maize is at least 3X that of ancestral teosinte, while our singleton-based estimate from the the SeeDs project data (citation) implies an  $N_e$  of close to 1 million individuals, or  $\sim 8X$  that of ancestral teosinte. From the perspective of a new mutation, both of these estimates are biased downward, since they represent past reductions in diversity due to the aforementioned bottleneck. A back-of-the-envelope calculation of modern  $N_e$  can be made by assuming there are 47.9 million ha of landrace maize in production [45], 42,000 plants per ha for open pollinated maize varieties [46], and conservatively that no more than one in 1,000 plants contribute gametes to the next generation. This calculation implies a modern  $N_e \geq 2$  billion indi-

viduals, and incorporates no individuals from hybrid breeding programs. If for some reason this back-of-the-envelope calculation is an overestimate, even by several orders of magnitude, it is clear that the post-domestication expansion in maize  $N_e$  is enormous. Patterns of singleton diversity in maize, as discussed in more detail in the next section, demonstrate that this recent expansion has a notable impact on the maize genome.

### Demography strongly impacts efficiency of purifying selection.

The observation that maize pairwise diversity is less impacted by the distance to the nearest gene than is teosinte pairwise diversity is reasonable from a long-term evolutionary standpoint. Theory has established that purifying selection is more efficient in a large population than in a small one [47], and this observation most likely reflects that prediction. More specifically, if teosinte  $N_e$  remained relatively constant while maize bottlenecked and recovered exponentially, this provides that the average  $N_e$  of maize over the previous several thousand generations is much smaller than that of teosinte, regardless of how much maize has ballooned in the recent past. Therefore, our observation shows that purifying selection in maize has not purged deleterious alleles, or the neutral alleles they are linked to, as effectively as in teosinte.

The reversal of this trend when we analyze only singleton diversity instead of pairwise diversity, however, stands out as a notable observation. Every mutation begins as a singleton (an allele present in only one individual), and therefore singletons are, on average, the youngest class of alleles that can be observed. Therefore unlike pairwise patterns of diversity, which are most heavily influenced by intermediate frequency alleles based on the definition of  $\pi$  [48], singleton diversity is most influenced by recent patterns of evolution. Hence, because our demographic estimation indicates dramatic expansion of maize  $N_e$  in the recent past, we expect for purifying selection to presently operate more efficiently in maize than in teosinte. This observation is very clearly demonstrated by the fact that singleton diversity in maize is more impacted by the distance to the nearest gene than it is in teosinte, as was shown in Figure 4.

A consequence of the inefficient purifying selection that maize experienced during its bottleneck is likely that it harbors more deleterious alleles segregating at intermediate frequency than does teosinte. This could be a part of the explanation of why maize inbreds have continued to improve over the past several decades [49]; if deleterious alleles tend to be recessive and are particularly frequent, they will have ample opportunities to display their phenotypes in inbreds. Our results also demonstrate that recent purifying selection in maize has become much more effective, potentially explaining the ongoing improvement of these inbreds as maize lines are continuously select

Ultimately, we have shown that purifying selection in maize has operated very differently than purifying selection in teosinte. This observation, along with the potential for soft sweeps, appear to explain the phenotypic divergence pretty hard to argue that purifying selection explains phenotypic divergence! between maize and teosinte much more completely than does positive selection generating hard-sweeps at protein-coding mutations. Importantly, our estimation of the parameters of the maize domestication bottleneck contribute to the understanding of how the demography of crop domestication can impact crop diversity. The bottleneck-effects from a sudden collapse in population size have been well studied and are known to impact crops for thousands of generations. Complementing this knowledge, our results demonstrate that the rapid expansion experienced by many crops after domestication can also have a profound influence on patterns of diversity, and the effects

of this expansion should be accounted for as important contributors to long-term evolution.

## Materials and Methods

**BASH, R, and Python scripts.** All scripts used for analysis are available in an online repository at [REPO ADDRESS HERE](#).

**Plant materials.** We made use of published sequence of inbred accessions of teosinte (*Z. mays* ssp. *parviglumis*) and maize landraces from the Maize HapMap2 panel [22,50]. We removed **X**teosinte and **X**maize individuals that appeared as outliers in an initial principal component analysis (Figure S5), leaving 13 teosinte and 23 maize that were used for all subsequent analyses (Table S6).

**Physical and genetic maps.** Sequences were mapped to the maize B73 version 3 reference genome [51] ([ftp://ftp.ensemblgenomes.org/pub/plants/release-22/fasta/zea\\_mays/dna/](ftp://ftp.ensemblgenomes.org/pub/plants/release-22/fasta/zea_mays/dna/)) as described by [22]. **did we use their bam's? I thought we used bwa bam's, which are not listed and weren't used for hapmap2. if that's the case, we need to mention this difference and acknowledge Robert for providing those bam's. either way, we need to include a link to the bamfiles too. somewhere here or above mention URL for Tripsacum bam's. did we align those too or download someone else's?** All analyses made use of uniquely mapping reads with mapping quality score  $\geq 30$  and bases with base quality score  $\geq 20$ ; quality scores around indels were adjusted following [52]. We converted physical coordinates to genetic coordinates via linear interpolation of the 1cM resolution NAM genetic map [53]. **i thought this should be wallace not glaubit's? need to specify what you do with physical positions before and after last marker on each chrom and markers, if any, that have physical and genetic order that disagree.**

**Estimating the site frequency spectrum.** We estimated both the genome-wide site frequency spectrum (SFS) as well as a separate SFS for genic (within annotated transcript) and intergenic ( $\geq 5kb$  from a transcript) regions. We used the biomaRt package [54,55] of R [56] to parse annotations from the **X**annotation **include annotation version number and link to ENSEMBL annotation used.** We estimated single population and joint SFS with the software ANGSD [57], including all positions with at least one aligned read in  $\geq 80\%$  of samples in one or both populations. We assumed individuals were fully inbred and treated each line as a single haplotype. **i removed options details because they don't make sense to casual reader and the interested reader can find the script.** Because ANGSD cannot calculate a folded joint SFS, we first polarized SNPs using the maize reference genome and then folded spectra using  $\delta\alpha\delta i$  [9].

**Demographic inference.** We used the software  $\delta\alpha\delta i$  [9] to estimate parameters of a domestication bottleneck from the joint maize-teosinte SFS. We modeled a teosinte

population of constant effective size  $N_a$ , that at time  $T_b$  generations in the past gives rise to a maize population of size  $N_b$  which grows exponentially to size  $N_m$  in the present (Figure 3) The model includes migration of  $M_{mt}$  individuals each generation from maize to teosinte and  $M_{tm}$  individuals from teosinte to maize. We fixed  $N_a$  using  $\delta\alpha\delta i$ 's estimation of  $\theta = 4N_a\mu$  from the data and a mutation rate of  $3 \times 10^{-8}$  [33]. We estimated all other parameters using 1,000  $\delta\alpha\delta i$  optimizations and allowing initial values to be randomly perturbed by a factor of 2. Initial values and upper and lower bounds can be found in table **X**. **this looked terrible in text, please make a supp. table. keep track of notation as e.g. we explain  $N_b$  as population size of bneck, but in reality parameter we fit is  $\frac{N_b}{N_a}$ .** We report parameter estimates from the optimization run with the highest log-likelihood.

We further made use of a large genotyping data set of more than 4,000 maize landraces [58] to estimate the modern maize  $N_e$  from singleton counts. We filtered these data to include only SNPs with data in  $\geq 1,500$  individuals, and then projected the SFS down to a sample of 500 individuals. **cite or mention how projection was done.** We then estimated  $N_e$  from the data assuming  $\mu = 3 \times 10^{-8}$  [33] and the relation  $4N_e\mu = \frac{S}{L}$  [34], where  $S$  is the total number of singleton SNPs and  $L$  is the total number of SNPs in the dataset.

Finally, [MSMC text here](#)

**Diversity.** We made use of the software ANGSD [57] for diversity calculations and genotype calling. We calculated diversity statistics in maize and teosinte in 100bp non-overlapping windows using filters as described above for the SFS. **this said 1,000 kb which is 1Mb, but should say 100bp right?** We used allele counts to estimate the number of singleton polymorphisms in each window, and used binomial sampling to create a second maize data set down-sampled to have the same number of samples as teosinte. We called genotypes in maize, teosinte, and *Tripsacum* at sites with a SNP p-value  $< 10^{-6}$  and when the genotype posterior probability  $> 0.95$ . We identified substitutions in maize and teosinte as all sites with a fixed difference with *Tripsacum* and  $\leq 20\%$  missing data. Substitutions were classified as synonymous, missense, or noncoding using the ensembl variant effects predictor [59]. For each window with  $\geq 100bp$  of data we computed the genetic distance between the window center and the nearest synonymous and missense substitution as well as the genetic distance to the nearest gene transcript. **this said gene center, but we're not using gene center anymore, right?**

**Selection scan.** We scanned the genome to identify sites that have experienced recent positive selection using the H12 statistic [37] in sliding windows of 200 SNPs with a step of 25 SNPs.

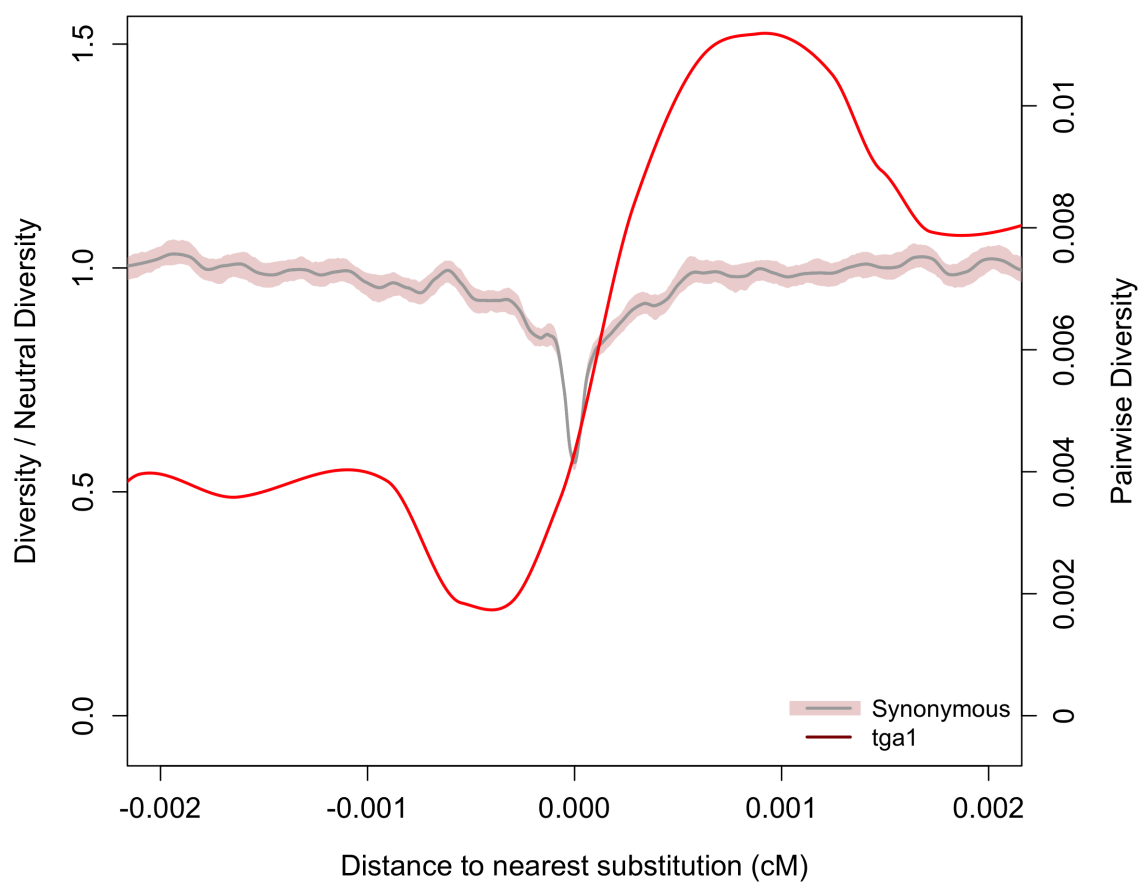
**ACKNOWLEDGMENTS.** We are indebted to Graham Coop and Simon Aeshbacher for their constructive input during this study. Funding was provided by NSF Plant Genome Research Project 1238014.

- Doebley, J. F., Gaut, B. S., & Smith, B. D. (2006) Cell 127, 1309–1321.
- Purugganan, M. D. & Fuller, D. Q. (2009) Nature 457, 843–848.
- Ross-Ibarra, J., Morrell, P. L., & Gaut, B. S. (2007) Proceedings of the National Academy of Sciences 104, 8641–8648.
- Hufford, M. B., Xu, X., Van Heerwaarden, J., Pyhäjärvi, T., Chia, J.-M., Cartwright, R. A., Elshire, R. J., Glaubitz, J. C., Guill, K. E., Kaeppeler, S. M., et al. (2012) Nature genetics 44, 808–811.
- He, Z., Zhai, W., Wen, H., Tang, T., Wang, Y., Lu, X., Greenberg, A. J., Hudson, R. R., Wu, C.-I., & Shi, S. (2011) PLoS genetics 7, e1002100.
- Vigouroux, Y., McMullen, M., Hittinger, C., Houchins, K., Schulz, L., Kresovich, S., Matsuoka, Y., & Doebley, J. (2002) Proceedings of the National Academy of Sciences 99, 9650–9655.
- Chapman, M. A., Pashley, C. H., Wenzler, J., Hvala, J., Tang, S., Knapp, S. J., & Burke, J. M. (2008) The Plant Cell 20, 2931–2945.
- Luijck, G., England, P. R., Tallmon, D., Jordan, S., & Taberlet, P. (2003) Nature Reviews Genetics 4, 981–994.
- Gutenkunst, R. N., Hernandez, R. D., Williamson, S. H., & Bustamante, C. D. (2009) PLoS genetics 5, e1000695.
- Li, J., Li, H., Jakobsson, M., Li, S., Sjödin, P., & Lascoux, M. (2012) Molecular Ecology 21, 28–44.
- Slotte, T. (2014) Briefings in functional genomics 13, 268–275.
- Sella, G., Petrov, D. A., Przeworski, M., & Andolfatto, P. (2009) PLoS genetics 5, e1000495.
- Gazave, E., Ma, L., Chang, D., Coventry, A., Gao, F., Muzny, D., Boerwinkle, E., Gibbs, R. A., Sing, C. F., Clark, A. G., et al. (2014) Proceedings of the National Academy of Sciences 111, 757–762.
- Shiferaw, B., Prasanna, B. M., Hellin, J., & Bänziger, M. (2011) Food Security 3, 307–327.
- Strable, J. & Scanlon, M. J. (2009) Cold Spring Harbor Protocols 2009, pdb-emo132.
- Smith, B. D. (1995) The emergence of agriculture. (Scientific American Library New York).
- Matsuoka, Y., Vigouroux, Y., Goodman, M. M., Sanchez, J., Buckler, E., & Doebley, J. (2002) Proceedings of the National Academy of Sciences 99, 6080–6084.
- Wilkes, H. G. et al. (1967) Teosinte: the closest relative of maize.
- Doebley, J., Stec, A., & Gustus, C. (1995) Genetics 141, 333.
- Wills, D. M., Whipple, C. J., Takuno, S., Kursel, L. E., Shannon, L. M., Ross-Ibarra, J., & Doebley, J. F. (2013) PLoS Genet 9, e1003604.
- Wang, H., Studer, A. J., Zhao, Q., Meeley, R., & Doebley, J. F. (2015) Genetics pp. genetics–115.
- Chia, J.-M., Song, C., Bradbury, P. J., Costich, D., de Leon, N., Doebley, J., Elshire, R. J., Gaut, B., Geller, L., Glaubitz, J. C., et al. (2012) Nature genetics 44, 803–807.
- Wright, S. I., Bi, I. V., Schroeder, S. G., Yamasaki, M., Doebley, J. F., McMullen, M. D., & Gaut, B. S. (2005) Science 308, 1310–1314.
- Wang, R.-L., Stec, A., Hey, J., Lukens, L., & Doebley, J. (1999) Nature 398, 236–239.
- Eyre-Walker, A., Gaut, R. L., Hilton, H., Feldman, D. L., & Gaut, B. S. (1998) Proceedings of the National Academy of Sciences 95, 4441–4446.
- Ross-Ibarra, J., Tenailon, M., & Gaut, B. S. (2009) Genetics 181, 1399–1413.
- Enard, D., Messer, P. W., & Petrov, D. A. (2014) Genome research 24, 885–895.
- Davydov, E. V., Goode, D. L., Sirota, M., Cooper, G. M., Sidow, A., & Batzoglou, S. (2010) PLoS Comput Biol 6, e1001025.
- Cooper, G. M., Stone, E. A., Asimenos, G., Green, E. D., Batzoglou, S., & Sidow, A. (2005) Genome research 15, 901–913.

30. Rodgers-Melnick, E, Bradbury, P. J, Elshire, R. J, Glaubitz, J. C, Acharya, C. B, Mitchell, S. E, Li, C, Li, Y, & Buckler, E. S. (2015) *Proceedings of the National Academy of Sciences* 112, 3823–3828.
31. Charlesworth, B, Morgan, M, & Charlesworth, D. (1993) *Genetics* 134, 1289–1303.
32. Tajima, F. (1989) *Genetics* 123, 585–595.
33. Clark, R. M, Tavaré, S, & Doebley, J. (2005) *Molecular biology and evolution* 22, 2304–2312.
34. Fu, Y.-X & Li, W.-H. (1993) *Genetics* 133, 693–709.
35. Wilson, B. A, Petrov, D. A, & Messer, P. W. (2014) *Genetics* 198, 669–684.
36. Innan, H & Kim, Y. (2004) *Proceedings of the National Academy of Sciences of the United States of America* 101, 10667–10672.
37. Garud, N. R, Messer, P. W, Buzbas, E. O, & Petrov, D. A. (2015) *PLoS genetics* 11, e1005004.
38. Clark, R. M, Wagler, T. N, Quijada, P, & Doebley, J. (2006) *Nature genetics* 38, 594–597.
39. Studer, A, Zhao, Q, Ross-Ibarra, J, & Doebley, J. (2011) *Nature genetics* 43, 1160–1163.
40. Beissinger, T. M, Hirsch, C. N, Vaillancourt, B, Deshpande, S, Barry, K, Buell, C. R, Kaeppler, S. M, Gianola, D, & de Leon, N. (2014) *Genetics* 196, 829–840.
41. Messer, P. W & Petrov, D. A. (2013) *Trends in ecology & evolution* 28, 659–669.
42. Sattath, S, Elyashiv, E, Kolodny, O, Rinott, Y, & Sella, G. (2011) *PLoS genetics* 7, e1001302.
43. Williamson, R, Josephs, E, Platt, A, Hazzouri, K, Haudry, A, Blanchette, M, & Wright, S. (2014) *PLoS genetics* 10, e1004622–e1004622.
44. Hernandez, R. D, Kelley, J. L, Elyashiv, E, Melton, S. C, Auton, A, McVean, G, Sella, G, Przeworski, M, et al. (2011) *science* 331, 920–924.
45. Program, T. M. (1999) *Development, maintenance, and seed multiplication of open-pollinated maize varieties. (CIMMYT, Mexico, D.F.), 2 edition.*
46. Van Heerwaarden, J, Van Eeuwijk, F, & Ross-Ibarra, J. (2010) *Heredity* 104, 28–39.
47. Kimura, M. (1984) *The neutral theory of molecular evolution. (Cambridge University Press).*
48. Nei, M & Li, W.-H. (1979) *Proceedings of the National Academy of Sciences* 76, 5269–5273.
49. Meghji, M, Dudley, J, Lambert, R, & Sprague, G. (1984) *Crop Science* 24, 545–549.
50. Lemmon, Z. H, Bukowski, R, Sun, Q, & Doebley, J. F. (2014) *PLoS Genet* 10, e1004745.
51. Schnable, P. S, Ware, D, Fulton, R. S, Stein, J. C, Wei, F, Pasternak, S, Liang, C, Zhang, J, Fulton, L, Graves, T. A, et al. (2009) *science* 326, 1112–1115.
52. Li, H. (2011) *Bioinformatics* 27, 2987–2993.
53. Glaubitz, J. C, Casstevens, T. M, Lu, F, Harriman, J, Elshire, R. J, Sun, Q, & Buckler, E. S. (2014) *PLoS One* 9, E90346.
54. Durinck, S, Spellman, P. T, Birney, E, & Huber, W. (2009) *Nature protocols* 4, 1184–1191.
55. Durinck, S, Moreau, Y, Kasprzyk, A, Davis, S, De Moor, B, Brazma, A, & Huber, W. (2005) *Bioinformatics* 21, 3439–3440.
56. R Core Team. (2014) *R: A Language and Environment for Statistical Computing (R Foundation for Statistical Computing, Vienna, Austria).*
57. Korneliussen, T. S, Albrechtsen, A, & Nielsen, R. (2014) *BMC bioinformatics* 15, 356.
58. Hearne, S, Chen, C, Buckler, E, & Mitchell, S. (2015) *Unimputed gbs derived snps for maize landrace accessions represented in the seed-maize gwas panel (). Accessed: 2015-02-16.*
59. McLaren, W, Pritchard, B, Rios, D, Chen, Y, Flicek, P, & Cunningham, F. (2010) *Bioinformatics* 26, 2069–2070.



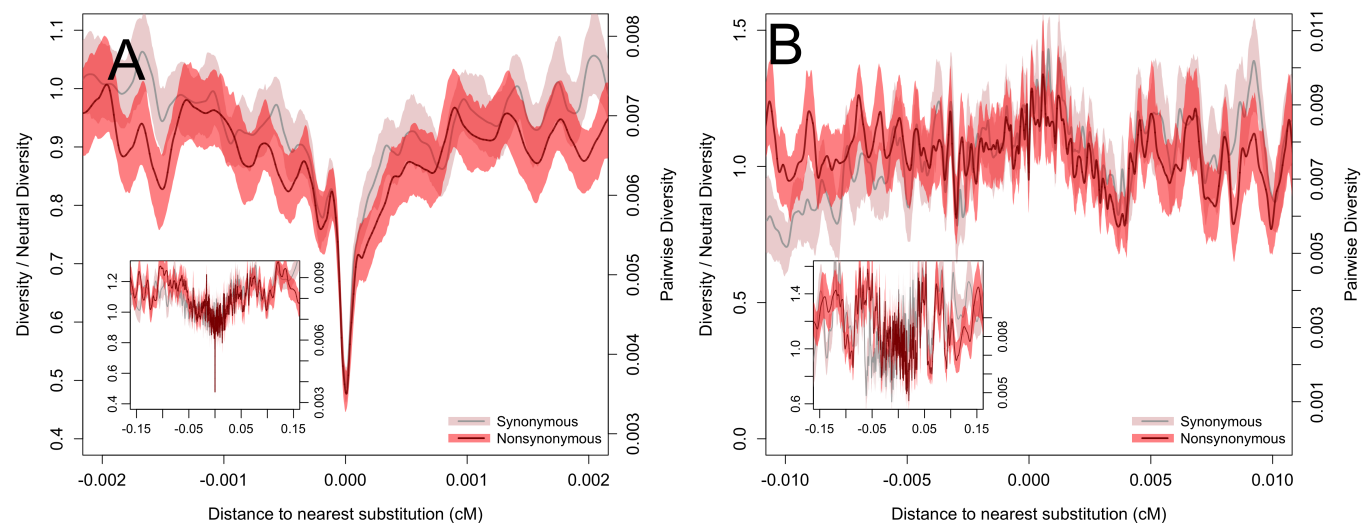




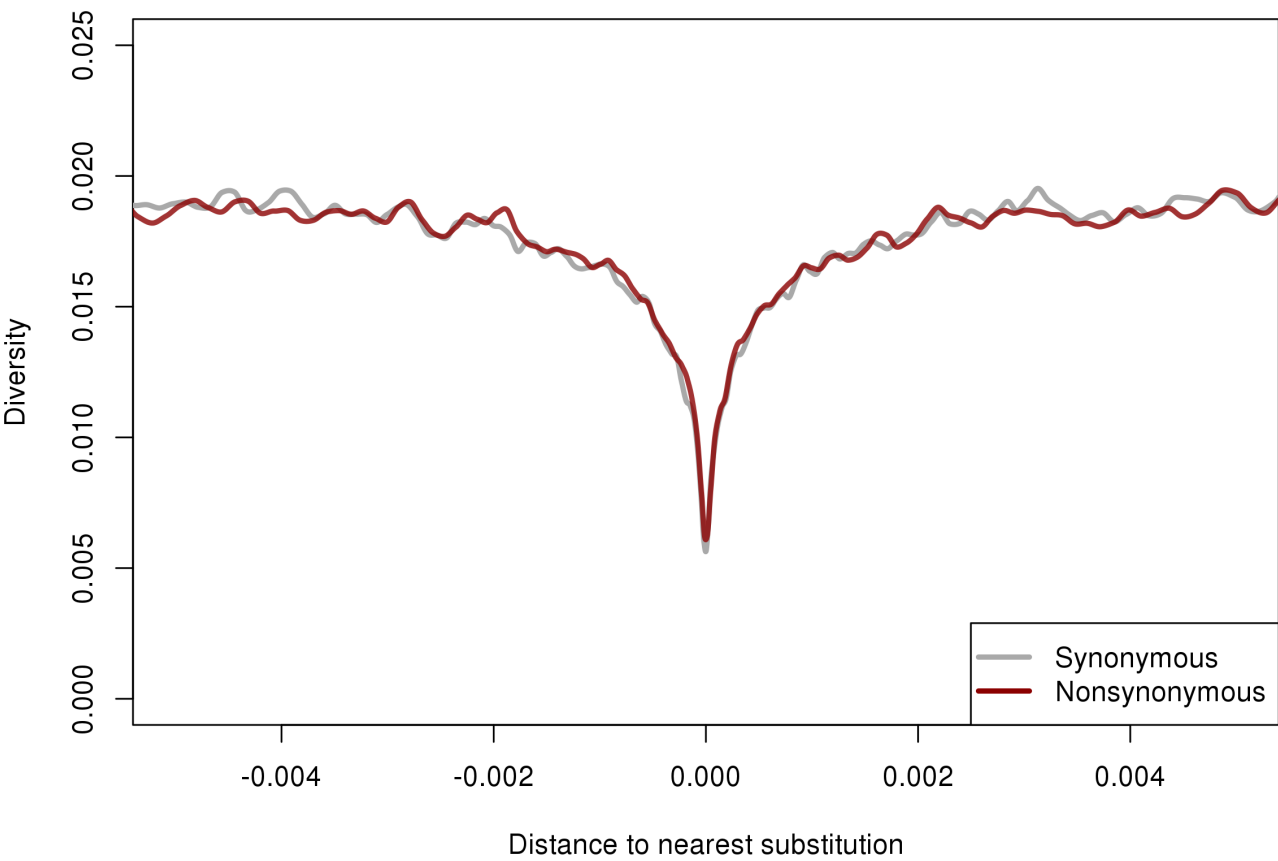
**Fig. S1.** Diversity surrounding the causative polymorphism at the *tga1* locus is plotted. Since this is only one gene, the large amount of noise compared to our average plots is expected. However, notice that diversity precisely at the causative polymorphism is reduced and a recovery of diversity is observed away from that site.

## Supporting Information

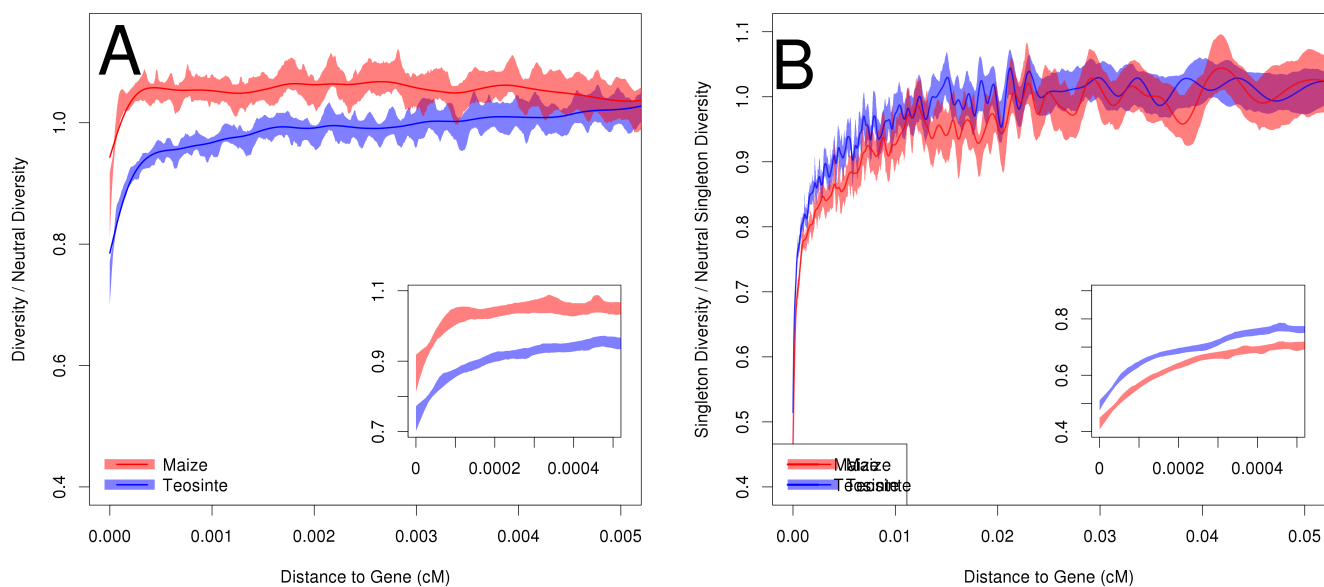




**Fig. S2.** Pairwise diversity surrounding synonymous and nonsynonymous substitutions in maize at highly conserved (A) or unconserved (B) sites. Bootstrap-based 95% confidence intervals are depicted via shading. Inset plots depict a larger range on the x-axis.

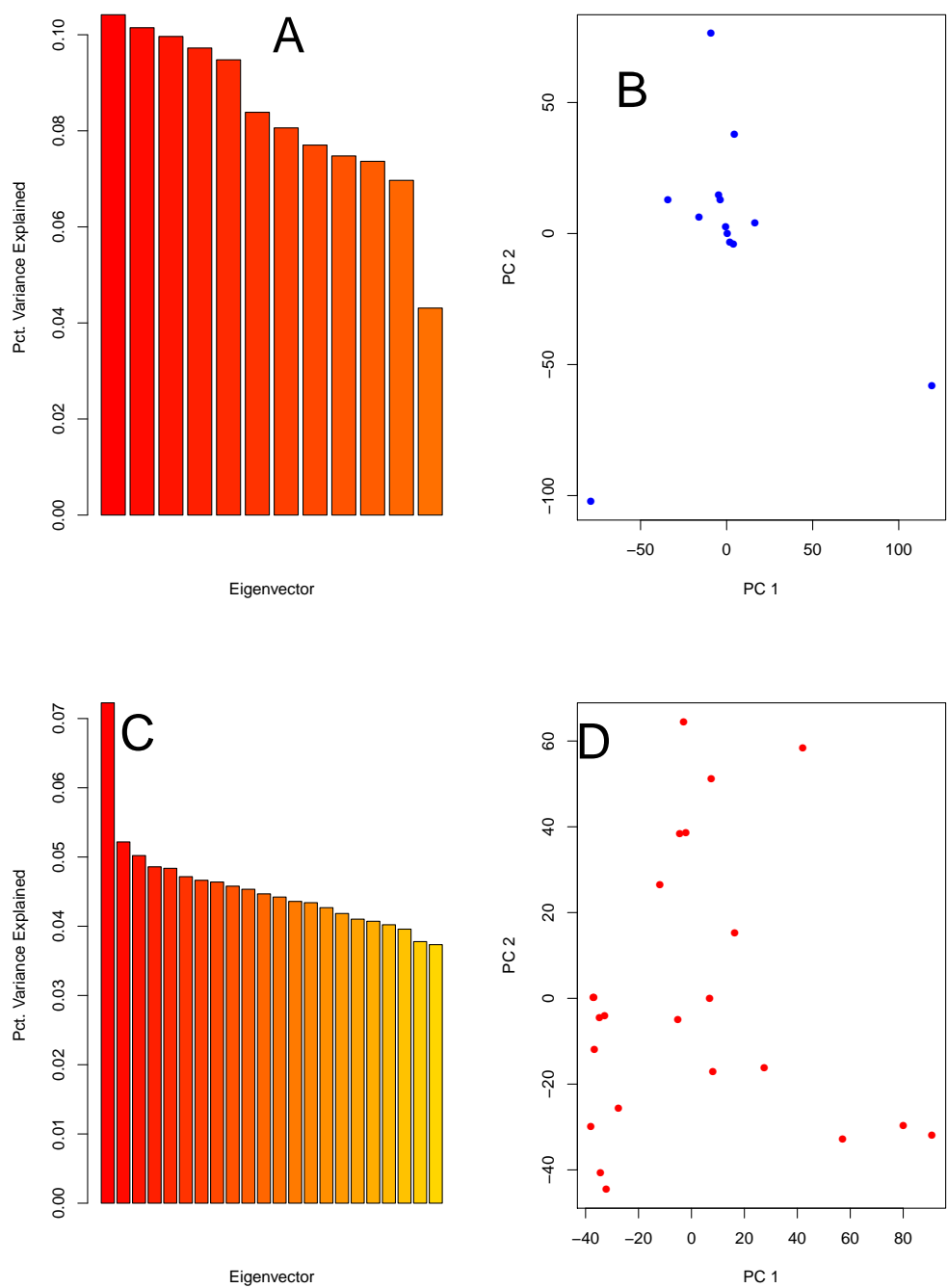


**Fig. S3.** Singleton diversity surrounding synonymous and nonsynonymous substitutions in maize.



**Fig. S4.** Relative level of diversity versus distance to the nearest gene, in maize and teosinte, based on only sites that do not show evidence of hard or soft sweeps according to H12. Two measures of diversity were investigated. **A** displays pairwise diversity, which is most influenced by intermediate frequency alleles and therefore depicts more ancient evolutionary patterns, and **B** depicts singleton diversity, influenced by rare alleles and thus depicting evolutionary patterns in the recent past. Bootstrap-based 95% confidence intervals are depicted via shading. Inset plots depict a smaller range on the x-axis.





**Fig. S5.** Principal component analysis of teosinte and maize individuals to ensure that no close relatives were inadvertently included in our study. Plots are based on a random sample of 10,000 SNPs. **A:** Percentage of total variance explained by each principal component for teosinte. **B:** PC1 vs PC2 for all 13 teosinte individuals. **C:** Percentage of total variance explained by each principal component for maize. **D:** PC1 vs PC2 for all 23 maize individuals. need to fix letter label positions

Maize	Teosinte
BKN009	TIL01
BKN010	TIL02
BKN011	TIL03
BKN014	TIL04-TIP454
BKN015	TIL07
BKN016	TIL09
BKN017	TIL10
BKN018	TIL11
BKN019	TIL12
BKN020	TIL14-TIP498
BKN022	TIL15
BKN023	TIL16
BKN025	TIL17
BKN026	
BKN027	
BKN029	
BKN030	
BKN031	
BKN032	
BKN033	
BKN034	
BKN035	
BKN040	

**Fig. S6.** A list of maize and teosinte individuals included in this study. Sequencing and details were previously described by [cite chia and lemmon](#)