

# Historical Divergence and Gene Flow in the Genus *Zea*

Jeffrey Ross-Ibarra,<sup>\*,1</sup> Maud Tenaillon<sup>†</sup> and Brandon S. Gaut<sup>\*,2</sup>

<sup>\*</sup>Department of Ecology and Evolutionary Biology, University of California, Irvine, California 92697 and <sup>†</sup>CNRS, UMR 0320/UMR 8120 Génétique Végétale, F-91190 Gif sur Yvette, France

Manuscript received October 16, 2008  
Accepted for publication January 8, 2009

## ABSTRACT

Gene flow plays a fundamental role in plant evolutionary history, yet its role in population divergence—and ultimately speciation—remains poorly understood. We investigated gene flow and the modalities of divergence in the domesticated *Zea mays* ssp. *mays* and three wild *Zea* taxa using sequence polymorphism data from 26 nuclear loci. We described diversity across loci and assessed evidence for adaptive and purifying selection at nonsynonymous sites. For each of three divergence events in the history of these taxa, we used approximate Bayesian simulation to estimate population sizes and divergence times and explicitly compare among alternative models of divergence. Our estimates of divergence times are surprisingly consistent with previous data from other markers and suggest rapid diversification of lineages within *Zea* in the last ~150,000 years. We found widespread evidence of historical gene flow, including evidence for divergence in the face of gene flow. We speculate that cultivated maize may serve as a bridge for gene flow among otherwise allopatric wild taxa.

GENE flow plays a fundamental role in plant evolutionary history, from maintaining cohesion among geographically separated populations (MORJAN and RIESEBERG, 2004; ARNOLD 2006) to accelerating evolution via adaptive introgression (ARNOLD 2004). Yet despite its importance, the role of gene flow in population divergence and speciation is poorly understood. Conventional theory argues strongly for the predominance of an allopatric model of divergence, in which populations are separated geographically and evolve in isolation without genetic exchange (COYNE and ORR 2004). Nonetheless, the simple allopatric model may not be appropriate for many plant lineages: ample plant data support the occurrence of population divergence (ANTONOVICS 2006; ARNOLD 2006; MALLET 2007; RIESEBERG and WILLIS 2007) and perhaps even speciation (SAVOLAINEN *et al.* 2006) in the face of gene flow.

While the population genetics of gene flow have been studied in numerous plant species (HAMRICK and GODT 1989; HAMRICK and NASON 1996; COYNE and ORR 2004; MORJAN and RIESEBERG 2004), few studies have made explicit attempts to study the process of population divergence. In fact, conventional analyses of gene flow often make assumptions about drift–migration equilibrium that preclude simultaneous evaluation of diver-

gence and gene flow. Moreover, the vast majority of studies to date have not used DNA sequence data, instead relying on other molecular markers for which model-based inference is less tractable. Explicit molecular population genetic analysis of divergence between plant lineages has so far been limited to a handful of studies in *Arabidopsis* (RAMOS-ONSINS *et al.* 2004), rice (ZHANG and GE 2007), and tomatoes (STÄDLER *et al.* 2005, 2008). These studies have advanced our understanding of divergence beyond a simple analysis of gene flow under equilibrium conditions, using coalescent methods to estimate divergence parameters and test the null model of divergence in isolation. But while all of these studies have suggested that their data provide evidence of introgression, none have been able to statistically reject a null model of isolation or draw firm conclusions about the role of gene flow in plant speciation.

The fact that no clear picture has emerged from these initial studies stems from several important limitations of the work to date. First, investigations have been limited to retrospective assessment of divergence between heterospecific populations, and such analyses may paint a different picture than studies of populations currently undergoing both divergence and gene flow. Second, studies to date have looked at only a handful ( $\leq 10$ ) of loci, such that sampling variance may have as much effect as biological processes in determining observed patterns of variation. Third, and most importantly, these studies have all relied heavily on the divergence population genetic (DPG) approach of Hey and colleagues (WAKELEY and HEY 1997; WANG *et al.*

Sequence data from this article have been deposited with the EMBL/GenBank Data Libraries under accession nos. FJ707611–FJ708483.

<sup>1</sup>Present address: Department of Plant Sciences, University of California, One Shields Ave., Davis, CA 95616.

<sup>2</sup>Corresponding author: Department of Ecology and Evolutionary Biology, University of California, 5205 McCaugh Hall, Irvine, CA 92697.  
E-mail: bgaut@uci.edu

1997), which uses coalescent simulations to test the fit of observed data to a model of divergence in isolation. While these approaches represent a substantial advancement, the goodness-of-fit test employed in these analyses is likely to be conservative for identifying introgression (STÄDLER *et al.* 2008). More importantly, while newer software allows estimation of migration parameters (HEY and NIELSEN 2004; BECQUET and PRZEWSKI 2007), the DPG approach and its descendants do not evaluate alternative models of divergence (but see HEY and NIELSEN 2007 for an improved implementation).

Here we extend the DPG approach to investigate divergence and speciation in the genus *Zea* (Poaceae). The genus *Zea* is composed of four species distributed across Mexico and into Central America. *Zea* has historically been divided into two sections (DOEBLEY and IRTIS 1980; IRTIS and DOEBLEY 1980). Section *Luxuriantes* consists of the perennial *Zea diploperennis*, its autotetraploid derivative *Z. perennis*, and the annual *Z. luxurians*. The only species in section *Zea* is *Z. mays*, which is an annual composed of four subspecies: the domesticated maize (ssp. *mays*), its wild progenitor ssp. *parviglumis*, and the wild taxa ssp. *mexicana* and ssp. *huehuetenangensis*. A number of authors have utilized genetic data to elucidate relationships within *Zea* (DOEBLEY *et al.* 1984; BUCKLER and HOLTSFORD 1996; HANSON *et al.* 1996; HILTON and GAUT 1998; TIFFIN and GAUT 2001; FUKUNAGA *et al.* 2005; VIGOUROUX *et al.* 2005), but substantial uncertainty remains about the evolutionary history of the genus, due in part to the complicating effects of hybridization and introgression (WILKES 1977; DOEBLEY 1990; FUKUNAGA *et al.* 2005). And while considerable attention has been directed toward maize domestication (EYRE-WALKER *et al.* 1998; HILTON and GAUT 1998; MATSUOKA *et al.* 2002; TENAILLON *et al.* 2004; WRIGHT *et al.* 2005), the process of divergence among other lineages in *Zea* remains unclear. In particular, little is known about the role gene flow has played in diversification within *Zea*, even in the well-studied case of domestication.

In this study we combined existing sequence polymorphism data with new resequencing data to create a combined set of 26 loci sampled from *Z. luxurians* and three subspecies of *Z. mays*: ssp. *mays*, ssp. *parviglumis*, and ssp. *mexicana* (Figure 1). This sampling allowed us to study divergence at three distinct levels: recent domestication, subspecies differentiation, and speciation. We evaluated patterns of diversity in these four taxa, testing for evidence of recent introgression and investigating the history of selection using divergence data from the sister genus *Tripsacum*. We then utilized coalescent simulations under an approximate Bayesian framework to simulate four alternative models of divergence, explicitly evaluating the relative probability of our models and using the models to estimate effective population sizes and divergence times.

## METHODS

**Sampling:** We analyzed rangewide samples of *Z. luxurians* and three subspecies of *Z. mays*: the domesticated maize (ssp. *mays*), its wild progenitor ssp. *parviglumis*, and the wild ssp. *mexicana*. Basic passport information on the samples used can be found in supplemental Table S1, and geographic information about the wild taxa is presented in Figure 1.

**DNA sequencing:** Sequencing methods followed TENAILLON *et al.* (2001, 2004). Amplified PCR products were cloned into a pGem TA cloning vector, and a single clone (single allele per individual) was then sequenced in both directions with BigDye chemistries on ABI automated sequencers. We reamplified and recloned individuals showing evidence of singleton polymorphisms in the initial alignment, sequencing five clones per individual per locus to confirm the sequence of the original allele. Singleton polymorphisms that could not be confirmed in this manner were assumed to be due to PCR error and were discarded from the analysis. We supplemented these data with sequence data from previous publications (supplemental Table S2) and with outgroup sequence from *Tripsacum dactyloides* at most loci. In total, we sampled sequence data from rangewide collections of an average of ~13 individuals per taxon at 26 loci in four taxa of *Zea* (supplemental Table S2). All new sequences have been submitted to GenBank (accession nos. FJ707611–FJ708483).

**Sequence analysis:** Sequences were initially aligned using the software Geneious (Biomatters) and then manually adjusted. Twelve of these loci were known genes, and previous determinations of exon structure were used for all analyses of these genes. For the remaining loci, we first identified open reading frames (ORFs) using a combination of the ORF-finding software at NCBI and in Geneious and using blastn searches of mRNA databases in GenBank. Final decisions on coding regions were based on assessment of the translation of potential regions and homology of the translated product with proteins found in GenBank. We calculated a suite of standard diversity statistics from nongapped, biallelic sites, including pairwise nucleotide diversity  $\theta_{\pi}$ , WATTERSON's (1975) estimator  $\theta_W$ , TAJIMA's (1989) *D* statistic,  $F_{ST}$ , and HUDSON and KAPLAN's (1985) *R*<sub>min</sub>, using the analysis package of software from the libsequence C++ library (THORNTON 2003). We compared levels of silent diversity among taxa using the recursion equations of HUDSON (1990) to estimate the multilocus value of  $\theta$  ( $= 4N\mu$ , where *N* is population size and  $\mu$  is mutation rate) for loci with samples from all four taxa. We estimated the likelihood of values of the population recombination rate  $\rho$  ( $= 4Nc$ , where *c* is the recombination rate per locus per generation) over a grid from 0 to 100, using the composite-likelihood approach of McVEAN *et al.* (2002) implemented in the software LDHAT.

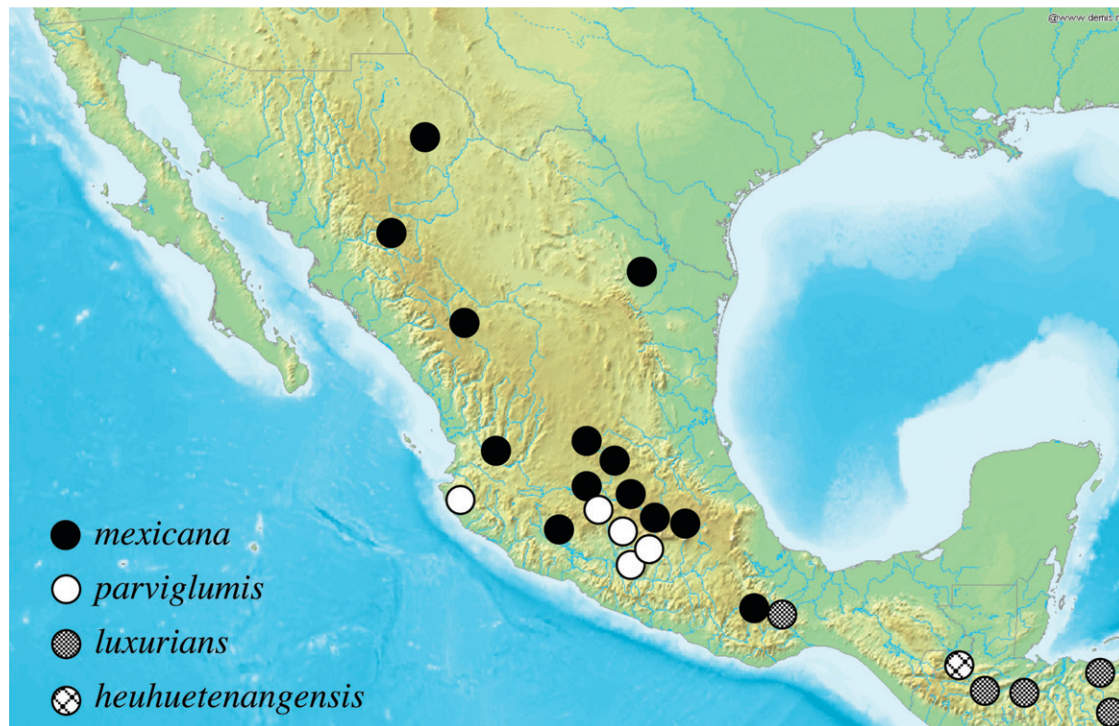


FIGURE 1.—Geographical distribution of *Zea luxurians* and the wild subspecies of *Z. mays*. Circles represent collection locations from herbarium specimens in the Missouri Botanical Garden and approximate the geographic range of these taxa.

**Detecting introgression:** We used the Bayesian clustering procedure of PRITCHARD *et al.* (2000) to test explicitly for recent introgression separately at each locus. We applied the program STRUCTURE (PRITCHARD *et al.* 2000) to data on the basis of two different genotype definitions. In the first, we defined haplotypes on the basis of sequence identity alone. In the second, we followed a procedure similar to CHEN *et al.* (2009), dividing each sequence into regions following the four-gamete test algorithm of HUDSON and KAPLAN (1985). Haplotypes within each region were defined by sequence polymorphism, excluding regions with fewer than three polymorphic sites. Thus, the genotype of each individual at a locus was defined ultimately by its haplotype configuration across regions. With each of the two genotypic definitions, we tested for gene flow among all four taxa, using the software STRUCTURE (PRITCHARD *et al.* 2000) to estimate the assignment probability of each individual to its original population. We ran STRUCTURE under the no admixture model for haploid data for 200,000 steps, using predefined populations assumed to have correlated allele frequencies, and with a burn-in of 50,000 steps including 10,000 steps of admix burn-in. Following suggestions in the documentation, we tested for evidence of admixture up to two generations in the past, assuming a prior probability of migration of 1%.

**Investigating selection:** We evaluated evidence for selection at these loci, using data on polymorphism and divergence from the 14 coding loci with samples from all four taxa and a *Tripsacum* outgroup (supplemental

Table S2). As an initial test for selection, we performed the McDonald–Kreitman (MK) test (MCDONALD and KREITMAN 1991) individually for all loci. We also performed the MK test on the pooled data, summing numbers of fixed and polymorphic synonymous and nonsynonymous sites across loci. For each taxon, we calculated the neutrality index (NI) (RAND and KANN 1996) as  $R_P S_F / R_F S_P$ , where  $R$  and  $S$  refer to counts of nonsynonymous and synonymous SNPs, and the  $P$  and  $F$  subscripts refer to polymorphic and fixed sites, respectively. An excess of fixed nonsynonymous sites relative to synonymous sites, often considered suggestive of positive selection, will lead to values of NI lower than unity, while an excess of polymorphic nonsynonymous sites may suggest purifying selection. Finally, we used counts of fixed and polymorphic sites to estimate the population selection parameter  $\gamma = 2Ns$ , where  $N$  is the effective population size and  $s$  is the selection coefficient. We estimated  $\gamma$  with the MKPRF software (BUSTAMANTE *et al.* 2002) for each locus separately as well as for data pooled across loci. We ran MKPRF with default parameters, except that burn-in and sampling were both extended to 5000 steps.

**Approximate Bayesian computation:** We used an approximate Bayesian framework (MARJORAM and TAVARE 2006) to investigate four models of population divergence. For each model, we compared *Z. mays* ssp. *parviglumis* against three other taxa: *Z. luxurians*, *Z. mays* ssp. *mexicana*, and domesticated maize. All four models are variations of a simple divergence model (WAKELEY and HEY 1997) in which an ancestral population in-



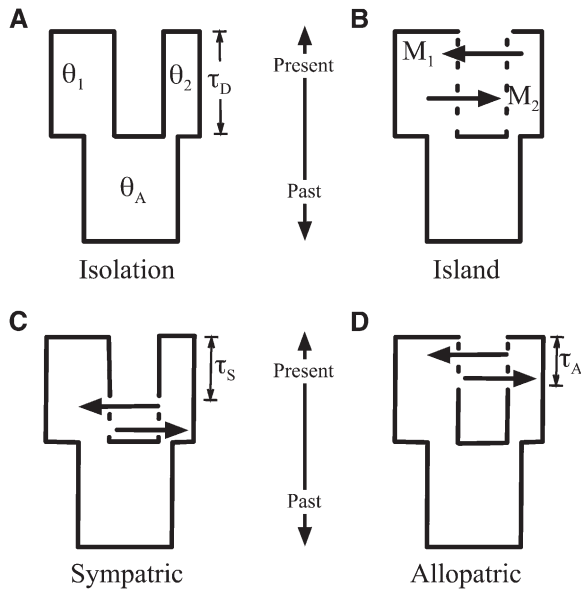


FIGURE 2.—Schematic representation of the four divergence models tested for each comparison. See text for description of parameters. (A) Isolation. (B) Island. (C) Sympatry. (D) Allopatry.

stantaneously splits into two daughter populations of constant size (Figure 2). Our first model, referred to here as the isolation model, assumes no gene flow between populations. The second model is an island model with continuous gene flow between populations. In the sympatric model, gene flow occurs only during population divergence and is followed by eventual isolation between populations. Our fourth model, termed the allopatric model, posits initial population divergence in isolation followed by renewed contact and gene flow.

All models incorporate the divergence time  $\tau_D$ , scaled in units of  $4N$  generations ( $N$  is the diploid effective population size), as well as the scaled population mutation rate  $\theta$  and the population recombination rate  $\rho$ . Mutation and recombination parameters are defined for the reference population ( $\theta_1$ ,  $\rho_1$ , referring to *Z. mays* ssp. *parviglumis* in each case), the second daughter population ( $\theta_2$  and  $\rho_2$ ), and the ancestral population ( $\theta_A$  and  $\rho_A$ ). Given the recent divergence between ssp. *parviglumis* and maize, we assumed  $\theta_A = \theta_1$  for this comparison (EYRE-WALKER *et al.* 1998; WRIGHT *et al.* 2005). For models with gene flow, migration is included as the scaled parameter  $M = 4Nm$  ( $M_1$  for migration into *Z. mays* ssp. *parviglumis* and  $M_2$  for migration into population 2), where  $m$  is the per generation probability of a lineage having immigrated from outside of the current population. Gene flow begins at a time  $\tau_A < \tau_D$  for the allopatric model and ceases at time  $\tau_S < \tau_D$  in the sympatric model. To convert scaled parameter estimates to values in terms of years or effective population size, we used a mutation rate of  $3 \times 10^{-8}$ /bp/year. This rate was estimated independently by CLARK *et al.* (2005) from

analysis of postdomestication mutations in locus *tb1* and calibrated using archaeological data on the timing of maize domestication. This rate is higher than previous estimates based on less certain fossil evidence (STEBBINS 1981; GAUT *et al.* 1996; WHITE and DOEBLEY 1999), but is nonetheless within the range of values thought plausible for plant nuclear sequences (WOLFE *et al.* 1987).

We simulated multilocus coalescent data under each of our models, drawing parameter values from specified prior distributions (supplemental Table S3). We used values of  $\theta_1$  and  $\rho_1$  estimated from the data, and, assuming that  $c$  and  $\mu$  do not vary among populations (but may vary among loci), placed uniform priors on the ratios  $\theta_A/\theta_1$  and  $\theta_2/\theta_1$ . Identical uniform prior distributions were also used for the divergence times  $\tau_D$ ,  $\tau_S$ , and  $\tau_A$ , although the latter two parameters were constrained as described above. Finally, the prior distribution of the migration parameters  $M_1$  and  $M_2$  was also assumed to follow a uniform distribution, but independent values of  $M_1$  and  $M_2$  were drawn for each locus.

To estimate parameter values from these simulations, we used Bayes' theorem. Specifically, we can write  $P(\Theta | X) \propto P(X | \Theta)P(\Theta)$ , where  $\Theta$  represents a vector of parameters of interest and  $X$  is the observed data. If the vector of summary statistics  $S$  is sufficient for the data, then  $P(\Theta | X) \propto P(S | \Theta)P(\Theta)$ . Given a proposed vector  $\Theta$ , we can simulate the summary statistic vector  $S'$ ; thus  $P(S | \Theta) = P(S' = S | \Theta)$ . To incorporate several continuously distributed summary statistics, we allow  $S'$  to approximate  $S$  by including a tolerance parameter  $\delta$ , such that  $P(S | \Theta) = P((S' - S)/S < \delta | \Theta)$ ; with a uniform prior probability  $P(\Theta)$  this becomes  $P(\Theta | X) \propto P((S' - S)/S < \delta | \Theta)$ , or the acceptance rate of simulations for a given vector  $\Theta$ .

For each simulated data set, we calculated the mean and variance across loci of four commonly used summary statistics (for a total of eight summary statistics): the number of shared ( $S_S$ ) and fixed ( $S_F$ ) silent sites between populations and the number of unique silent sites ( $S_{X1}$  and  $S_{X2}$ ) in each population (WAKELEY and HEY 1997; BECQUET and PRZEWSKI 2007; FAGUNDES *et al.* 2007). These summary statistics were compared to observed summary statistics from our sequence alignments, rejecting simulations not within  $\delta = 0.3$  of both the mean and the variance of each statistic. Posterior distributions for the parameters of interest were then estimated as described above. Software used for the generation of priors, rejection sampling, and coalescent simulation is available from the authors.

## RESULTS

**Diversity and divergence:** We collected sequence data at 26 loci from an average of 13 individuals from each of four *Zea* taxa (supplemental Tables S1 and S2). Silent diversity varies greatly among these taxa (Table 1), with maximum-likelihood estimates of the population muta-

TABLE 1

Diversity statistics at silent sites in the studied taxa of Zea

Taxon	<i>S</i>	$\eta_1$	<i>H</i>	$\theta_\pi$	<i>D</i>	RM	$\rho$
Maize	20.9	6.8	0.802	0.013	0.002	2.1	0.032
<i>parviglumis</i>	25.7	13.4	0.936	0.017	−0.402	2.7	0.059
<i>mexicana</i>	24.4	12.1	0.906	0.016	−0.295	2.5	0.063
<i>luxurians</i>	11.3	5.1	0.807	0.010	−0.082	1.0	0.051

Shown are average numbers of segregating sites (*S*) and singletons ( $\eta_1$ ), haplotype diversity (*H*), nucleotide diversity ( $\theta_\pi$ ), Tajima's *D* (*D*), Hudson and Kaplan's Rmin (RM), and the recombination rate per base pair ( $\rho$ ).

tion rate  $\theta$  in *Z. mays* ssp. *parviglumis* (hereafter simply "*parviglumis*") and *Z. mays* ssp. *mexicana* ("*mexicana*") significantly higher than estimates for *Z. mays* ssp. *mays* ("maize") and *Z. luxurians* ("*luxurians*") (Figure 3). Measures of nucleotide and haplotype diversity show a similar pattern (Table 1), pointing to a much larger effective population size in *parviglumis* and *mexicana*. Both of these taxa also contain an excess of low-frequency polymorphisms (high numbers of singletons and a negative Tajima's *D*). In contrast, maize and *luxurians* show a frequency spectrum shifted toward zero, with fewer singleton polymorphisms and Tajima's *D* values close to the expectation under neutral equilibrium. Differences in effective population size are reflected in our estimates of the population recombination rate  $\rho$  as well, with lower estimates for *luxurians* and maize (Table 1).

Substantial variation is evident among loci (Table 2 and supplemental Tables S4–S6). Within *parviglumis*, for example,  $\theta_\pi$  at silent sites ranges from 0.7% at locus *gl1510* to 3.6% at *d8* (Table 2). Loci also differ in frequency spectra, with Tajima's *D* values in *parviglumis* ranging from −1.68 at locus *d8* to 1.55 at *csu1138*. The *D* value of several loci deviates significantly from the expectation of the standard neutral model, but only *tb1* in maize remains significant after multiple-test correction (data not shown). In *parviglumis*, estimates of  $\rho$  per base pair are 0 for five loci (*bnl7-13*, *c1*, *d8*, *csu636*, and *fus6*), and as high as 0.48 for locus *mgs3020* (Table 2).

We made an initial assessment of differentiation between taxa by calculating  $F_{ST}$  and counts of shared ( $S_S$ ), fixed ( $S_F$ ), and unique ( $S_{X1}$  and  $S_{X2}$ ) polymorphisms between pairs of taxa (Figure 4).  $F_{ST}$  clearly reveals a distinction between *luxurians* and *Z. mays*. Among the subspecies of *Z. mays*, *parviglumis* has the lowest median  $F_{ST}$  to *luxurians*, at 0.19. In contrast, the subspecies of *Z. mays* have substantially lower pairwise  $F_{ST}$  values. Among these taxa, *parviglumis* has lowest  $F_{ST}$  when compared to *mexicana*, perhaps suggesting large long-term population sizes in both taxa. Maize has a lower  $F_{ST}$  (0.066) with *mexicana* than with its wild ancestor, *parviglumis* (0.079). The distributions of  $S_S$

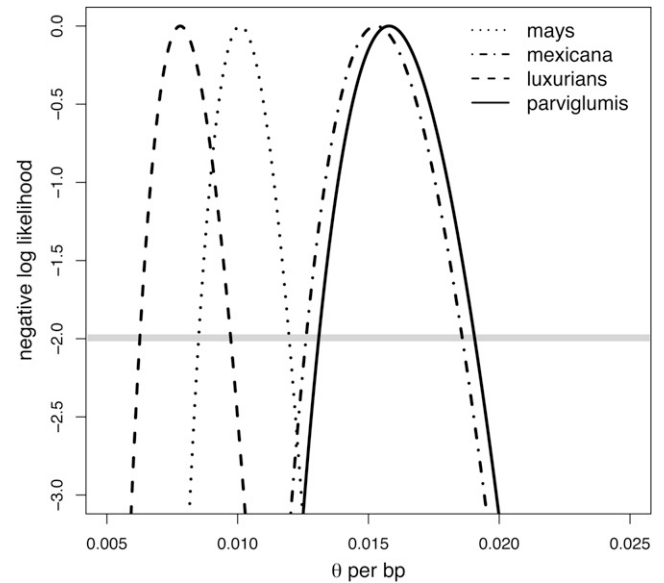


FIGURE 3.—Maximum-likelihood estimates of the population mutation rate  $\theta$  in four taxa of Zea. Ninety-five percent C.I.'s can be approximated by the intersection of each curve with the horizontal shaded bar at  $-2$ .

and  $S_F$  mirror the broad patterns revealed by  $F_{ST}$  (Figure 4).  $S_S$  is greater than zero in all pairwise comparisons among Zea taxa, but  $S_F > 0$  between *luxurians* and all subspecies of *Z. mays*, suggesting some measure of isolation between species. Among the subspecies of *Z. mays*, only *parviglumis* and maize have a fixed difference between them, due to one SNP at locus *asg65*.

**Selection:** If all mutations are neutral, drift alone should determine both the number of fixed differences between species and the level of polymorphism within species; the ratios of these values for different classes of SNPs are thus potentially informative about deviations from the expectation under neutrality (McDONALD and KREITMAN 1991). Using divergence information from a *Tripsacum* outgroup, we applied the MK test to the 14 loci with coding regions and with data representing each of our four Zea taxa (supplemental Table S2). None of the loci deviate significantly from neutral expectations in any taxon (data not shown). Consistent with the MK results, absolute values of the population selection parameter  $\gamma$  are low, with confidence intervals that overlap zero for most loci (Figure 5).

Analyses of individual loci may have low power to detect selection, and we thus analyzed samples pooled across loci. MK tests remain nonsignificant for counts pooled across loci, but the NI differs from unity for both *parviglumis* and maize (Table 3), suggesting deviation from neutrality in these taxa. For *parviglumis*, the NI of 0.76 differs from unity in a direction consistent with positive selection acting to fix nonsynonymous variants. For maize, an NI of 1.23 suggests the action of purifying selection preventing deleterious nonsynonymous variants from fixing. Estimates of  $\gamma$  from the pooled data

TABLE 2  
Diversity statistics at silent sites for all loci in *parviglumis*

Locus	bp	<i>n</i>	<i>S</i>	$\eta_1$	<i>H</i>	$\theta_\pi$	<i>D</i>	RM	$\rho$
<i>adh1</i>	808	8	57	27	0.964	0.026	-0.224	5	0.018
<i>asg11</i>	395	6	10	9	1.000	0.009	-1.069	0	0.001
<i>asg35</i>	467	12	31	19	0.985	0.016	-1.327	3	0.041
<i>asg64</i>	509	16	27	21	1.000	0.010	-1.505	1	0.058
<i>asg65</i>	531	14	20	4	0.912	0.013	0.297	2	0.005
<i>bnl7-13</i>	758	12	26	7	0.955	0.012	0.061	0	0.000
<i>bz2</i>	205	12	10	4	0.909	0.015	-0.255	1	0.037
<i>c1</i>	440	11	20	3	0.909	0.018	0.633	3	0.000
<i>csu1132</i>	330	8	28	16	0.964	0.029	-0.548	3	0.013
<i>csu1138</i>	186	14	5	1	0.835	0.012	1.547	1	0.047
<i>csu1171</i>	466	9	18	11	0.972	0.012	-0.869	1	0.010
<i>csu381</i>	916	10	42	24	1.000	0.014	-0.561	4	0.020
<i>csu636</i>	643	15	16	6	0.952	0.008	0.275	0	0.000
<i>csu838</i>	257	13	15	10	0.962	0.013	-1.186	2	0.120
<i>d8</i>	206	13	37	30	1.000	0.036	-1.677	2	0.000
<i>fus6</i>	153	10	4	2	0.644	0.010	0.143	0	0.000
<i>gl1510</i>	484	15	12	5	0.933	0.007	-0.437	2	0.137
<i>glb1</i>	511	8	52	37	1.000	0.033	-0.878	6	0.033
<i>mgs3020</i>	179	14	10	5	0.879	0.012	-1.254	2	0.483
<i>pepc1070</i>	712	13	63	22	1.000	0.029	0.098	8	0.016
<i>pepc1150</i>	233	16	15	8	0.883	0.015	-0.870	2	0.086
<i>tb1</i>	1982	7	79	49	0.952	0.015	-0.578	7	0.007
<i>ts2</i>	301	10	13	3	1.000	0.018	0.841	4	0.168
<i>vp1010</i>	167	12	9	4	0.803	0.017	-0.179	0	0.032
<i>wip1</i>	232	15	20	8	0.933	0.026	-0.088	4	0.060
<i>wx1</i>	497	19	29	14	0.994	0.013	-0.839	6	0.138

bp, number of silent sites; *n*, sample size; *S*, segregating sites;  $\eta_1$ , number of singletons; *H*, haplotype diversity;  $\theta_\pi$ , nucleotide diversity per base pair; *D*, Tajima's *D*; RM, minimum number of recombination events;  $\rho$ , recombination rate per base pair.

are consistent with NI estimates;  $\gamma$  is estimated to be close to zero for *mexicana* and *luxurians*, but is slightly positive in *parviglumis* and slightly negative in maize (Table 3).

The comparison of polymorphic to fixed sites highlights long-term evolutionary patterns. To investigate more recent evidence of selection, we plotted the frequency spectrum of synonymous and nonsynonymous polymorphisms pooled across all 14 loci (Figure 6). All four taxa show an excess of low-frequency nonsynonymous variants, a pattern commonly taken to represent the effects of purifying selection preventing weakly deleterious mutations from rising to high frequency or fixing in a population (Fu and Li 1993).

**Evidence for introgression:** Low  $F_{ST}$  values and shared variants suggest the potential for gene flow among *Zea* taxa, but these measures do not allow explicit tests of migration. We more formally tested for the possibility of gene flow at each locus using the software STRUCTURE (PRITCHARD *et al.* 2000), comparing a null model of no admixture (in which individuals are assumed to be assigned to the correct population) to an alternative model of recent admixture in the past several generations. These analyses are motivated by the observation of shared sequences

among taxa, including both entire sequences (supplemental Table S7) and sequence fragments. Using the entire sequence to define haplotypes, STRUCTURE analyses did not identify any likely cases of introgression (data not shown). Nonetheless, to examine the possibility that our observations of shared segments may be the result of past introgression events, we partitioned each locus into semi-independent fragments using the four-gamete test (see METHODS). Table 4 presents the results from these analyses, showing loci for which at least one genotype has a posterior probability of assignment to its original population of <50%. (*cf.* FUKUNAGA *et al.* 2005). Approximately half of the loci analyzed show evidence supporting introgression of sequence segments (Table 4). Most (~70%) of these inferences point to recent introgression among subspecies of *Z. mays*, although there is a suggestion of gene flow between *luxurians* and one or more subspecies of *Z. mays* at six loci.

**Modeling population divergence:** STRUCTURE analyses suggest there may have been recent gene flow among taxa. To further examine the historical role of gene flow during divergence, we employed an approximate Bayesian framework to assess the probability of four alternative models of population splitting (Figure

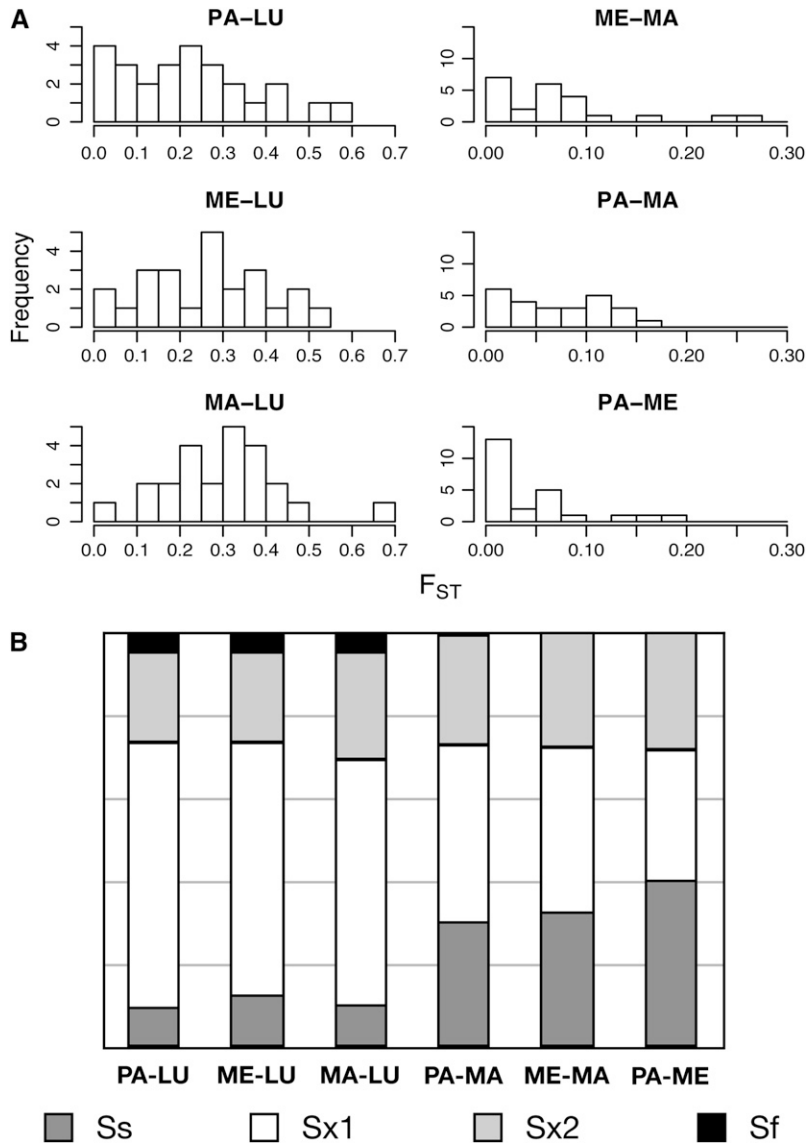


FIGURE 4.—Divergence in *Zea*. (A) Distribution of pairwise  $F_{ST}$  values between *Zea* taxa studied. (B) Comparison of proportions of shared ( $S_S$ ), fixed ( $S_F$ ), and unique ( $S_{X1}$  and  $S_{X2}$ ) silent sites among pairwise comparisons. MA, maize; ME, *mexicana*; PA, *parviglumis*; LU, *luxurians*.

2) in three independent divergence events. Our rejection-sampling scheme accepted simulated data that closely approximated the observed multilocus data. Using the acceptance rate  $P(S | Y_1)$ , or the probability of the observed vector of summary statistics  $S$  given model  $Y_1$ , we calculated a Bayes factor representing the relative weight of the evidence for alternative model  $Y_1$  over model  $Y_2$ , as  $K = P(S | Y_1) / P(S | Y_2)$ . Following JEFFREYS (1998), we considered values of  $K > 10^{1/2}$  “substantial” evidence for model  $Y_1$ , thus rejecting model  $Y_2$ . Bayes factors comparing each of the four models to the isolation model are shown in Table 5. Evaluation of the four models for the *parviglumis*–*luxurians* divergence allows us to reject alternative models in favor of the sympatric model, suggesting a history of ancestral gene flow during divergence. For the *parviglumis*–*mexicana* comparison, models excluding recent gene flow (isolation and sympatry) exhibit notably lower probabilities, but only the isolation model can be rejected.

The *parviglumis*–maize comparison reveals the opposite pattern: models including recent gene flow (island and allopatry) have the lowest probabilities, but only the island model can be rejected.

For each comparison, we performed parameter estimation under the most likely model by building a posterior density distribution of parameter values from accepted simulations (Figure 7). We took the highest probability value of each distribution as a point estimate of the parameter; using a recent estimate of the mutation rate in *Zea* (CLARK *et al.* 2005), we converted parameter estimates into values of effective population size ( $N$ ) and divergence time in years. While posterior probability distributions of the migration parameter  $M$  are not substantially different from our priors—suggesting that our data offer little power to quantitatively differentiate levels of gene flow among these taxa—distributions for all other parameters are clearly distinguishable. Estimates of effective population sizes



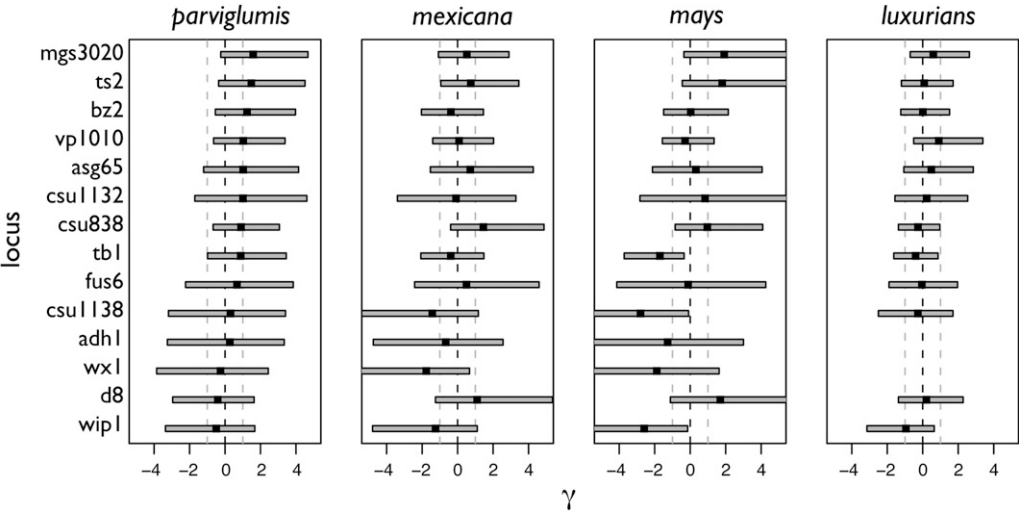


FIGURE 5.—Estimates of  $\gamma = 2N_s$ , based on MKPRF analyses. Solid squares represent point estimates, with shaded bars showing the 95% C.I. of each estimate. A black vertical dashed line represents  $\gamma = 0$  and gray dashed lines show the interval  $-1 < \gamma < 1$ . Loci *adh1* and *wxl* had too few sites in *luxurians* to estimate  $\gamma$ .

reveal relatively similar values of  $N$  across the phylogeny:  $N_s$  for *mexicana*, the *parviglumis*–*mexicana* ancestor, and the *parviglumis*–*luxurians* ancestor are all approximately the same as for modern-day *parviglumis*, or 120,000–160,000 individuals. The only estimates notably different from *parviglumis* are those of *luxurians* ( $\theta_2/\theta_1 \sim 40\%$  or 50,000 individuals) and maize ( $\theta_2/\theta_1 \sim 30\%$  or 45,000 individuals). Because our model restricts  $N$  to be constant within a lineage, however, these lower  $N$  values may represent strong population bottlenecks followed by recovery.

The posterior distributions of  $\tau_D$  all support divergence  $\ll 4N$  generations (Figure 7), consistent with the observed levels of shared polymorphism (Figure 4). Our estimate of  $\sim 55,000$  years for the *parviglumis*–maize divergence is nearly indistinguishable from the *mexicana* divergence estimate of  $\sim 60,000$  years, but severalfold higher than current dates for the domestication of maize (POHL *et al.* 2007); the estimated timing of isolation ( $\tau_S$ ) for *parviglumis*–maize is more recent ( $\sim 27,000$  years), but still implausible given the archeological record. Divergence ( $\tau_D$ ) between *parviglumis* and *luxurians* is estimated at  $\sim 140,000$  years and followed by the cessation of gene flow ( $\tau_A$ ) at  $\sim 60,000$  years in the past.

TABLE 3

Inference of selection in *Zea*

Taxa	SNP type	Fixed	Poly	NI	$\gamma$
Maize	Synonymous	119	90	1.23	−0.30
	Nonsynonymous	58	54		
<i>mexicana</i>	Synonymous	118	91	1.02	−0.06
	Nonsynonymous	61	48		
<i>luxurians</i>	Synonymous	123	64	1.08	−0.10
	Nonsynonymous	64	36		
<i>parviglumis</i>	Synonymous	108	135	0.76	0.33
	Nonsynonymous	62	59		

Poly, polymorphic; NI, neutrality index.

DISCUSSION

**Polymorphism and divergence:** We investigated diversity at 26 nuclear loci in four taxa of the genus *Zea* spanning three divergence events: between species, between subspecies, and between a domesticated and its progenitor. Even in the absence of coalescent models, our data provide information about the divergence history of the taxa studied. Given that the time to loss of shared variants is expected to be  $\sim 2N$  generations (CLARK 1997), the observed shared variation between *luxurians* and *Z. mays* (Figure 4) precludes an ancient divergence between these species without invoking either a very large effective population size for *luxurians* or a history of continual gene flow at all loci. Levels of recombination and diversity in *luxurians* (Table 1) attest to its reduced population size, however, and substantial unique polymorphism suggests that continual gene flow is unlikely. Among the subspecies of *Z. mays*, unique polymorphisms support the distinctness of *mexicana* and *parviglumis*, but patterns of shared variants and a low  $F_{ST}$  suggest a recent divergence and/or some amount of introgression (Figure 4).

The frequency spectrum of polymorphism provides information about changes in population size for these taxa. Negative values of Tajima’s  $D$  are not expected from specieswide sampling except under population size expansion (see discussion in ARUNYAWAT *et al.* 2007 and references therein). Along with previous evidence for population structure (DOEBLEY *et al.* 1984; FUKUNAGA *et al.* 2005; MOELLER *et al.* 2007), the observed negative  $D$  values in *parviglumis* and *mexicana* thus suggest a history of expansion, consistent with the fact that our coalescent estimate of  $N$  for the *parviglumis*–*mexicana* ancestor is roughly half that of the two descendant taxa combined. In contrast to *mexicana* and *parviglumis*, values of  $D$  near zero in maize and *luxurians* may reflect a reduction in effective population size. Multiple analyses of domestication have already suggested a population reduction for



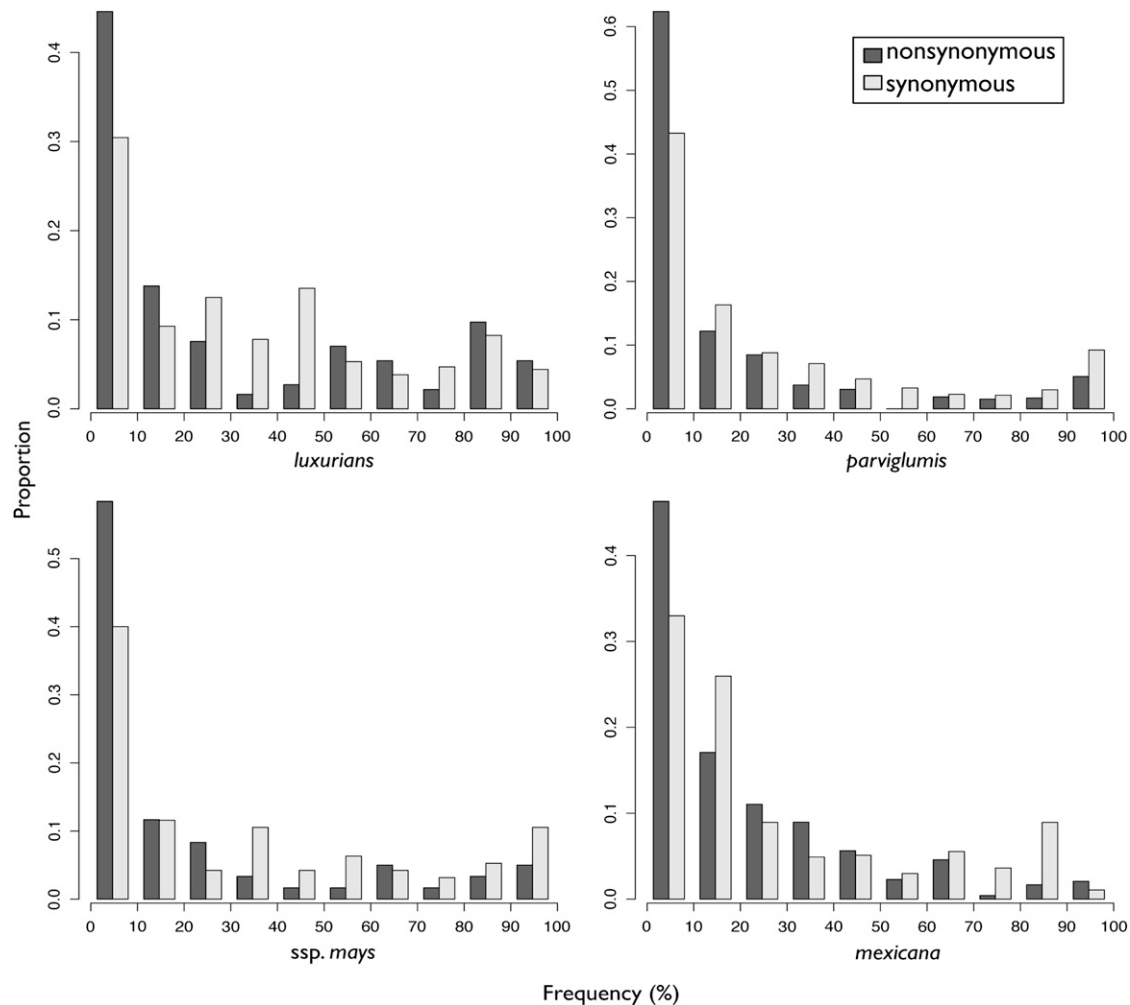


FIGURE 6.—Site frequency spectra of synonymous and nonsynonymous variants in each taxon. Frequency data have been binned into 10% intervals for plotting.

maize (EYRE-WALKER *et al.* 1998; HILTON and GAUT 1998; DOEBLEY 2004; TENAILLON *et al.* 2004; WRIGHT *et al.* 2005); for *luxurians*, the suggestion of a population reduction is consistent with previous studies (HILTON and GAUT 1998), low diversity levels (Figure 3), low  $\rho$ -estimates (Table 1), and a restricted geographic distribution (WILKES 1977; DOEBLEY and ILLIS 1980).

**Selection:** Few striking patterns emerge from our analyses of selection. In *luxurians*, confidence intervals for all estimates of the population selection parameter  $\gamma$  overlap zero (Figure 5), consistent with estimates of a low effective population size from our divergence model. A low  $N$  decreases the efficacy of selection, allowing otherwise strongly deleterious variants to rise to appreciable frequencies before being removed from the population and converting weakly deleterious mutations into effectively neutral variants. The site frequency spectrum of *luxurians* hints at this possibility; only in *luxurians* are nonsynonymous polymorphisms more frequent than synonymous polymorphisms for both low-frequency categories (0–20%; Figure 6).

The evidence for selection in the three subspecies of *Z. mays* is more difficult to interpret. Given their extensive shared evolutionary history, one might expect strong correlations in the patterns of selection that compare divergence from *Tripsacum* to extant polymorphism. Superficially, shared history might seem to explain the striking correlation in estimates of  $\gamma$  between *mexicana* and maize (Pearson's  $r > 0.85$ ; Figure 5), but the relatively distinct patterns of  $\gamma$  seen at loci in *parviglumis* (Pearson's  $r = 0.47$  and  $0.55$  for comparisons with *mexicana* and maize, respectively) cast doubt on such an interpretation. A potential explanation for the relatively low correlation of  $\gamma$  between *parviglumis* and *mexicana* is that differences in population structure obscure shared history (*e.g.*, FUKUNAGA *et al.* 2005). The striking correlation of  $\gamma$  between maize and *mexicana* could be indicative of recent gene flow (and thus shared history) between these taxa (BLANCAS *et al.* 2002; BALTAZAR *et al.* 2005; FUKUNAGA *et al.* 2005; ELLSTRAND *et al.* 2007); indeed, our analyses detect recent gene flow between maize and *mexicana* at several loci (Table 4).

**TABLE 4**  
**STRUCTURE analysis of putative admixture**

Locus	Taxon	
	Source	Destination
<i>adh1</i>	Maize	<i>parviglumis</i>
<i>adh1</i>	<i>mexicana</i>	<i>parviglumis</i>
<i>adh1</i>	<i>parviglumis</i>	Maize
<i>adh1</i>	Maize	<i>mexicana</i>
<i>asg35</i>	<i>parviglumis</i>	Maize
<i>c1</i>	Maize	<i>parviglumis</i>
<i>csu1132</i>	<i>mexicana</i>	<i>luxurians</i>
<i>csu1132</i>	<i>luxurians</i>	Maize
<i>csu1132</i>	<i>parviglumis</i>	Maize
<i>csu1132</i>	<i>luxurians</i>	<i>mexicana</i>
<i>csu1132</i>	Maize	<i>mexicana</i>
<i>csu1132</i>	Maize	<i>parviglumis</i>
<i>csu1138</i>	<i>parviglumis</i>	<i>luxurians</i>
<i>csu1171</i>	<i>parviglumis</i>	Maize
<i>csu381</i>	<i>luxurians</i>	Maize
<i>csu381</i>	<i>mexicana</i>	Maize
<i>csu381</i>	<i>luxurians</i>	<i>mexicana</i>
<i>csu381</i>	Maize	<i>mexicana</i>
<i>csu381</i>	Maize	<i>parviglumis</i>
<i>csu838</i>	<i>mexicana</i>	<i>luxurians</i>
<i>pepc1070</i>	<i>parviglumis</i>	Maize
<i>pepc1070</i>	<i>parviglumis</i>	<i>mexicana</i>
<i>pepc1070</i>	Maize	<i>parviglumis</i>
<i>pepc1150</i>	<i>mexicana</i>	Maize
<i>pepc1150</i>	<i>parviglumis</i>	<i>mexicana</i>
<i>pepc1150</i>	Maize	<i>mexicana</i>
<i>tb1</i>	<i>mexicana</i>	<i>luxurians</i>
<i>tb1</i>	Maize	<i>mexicana</i>
<i>tb1</i>	<i>luxurians</i>	<i>mexicana</i>
<i>tb1</i>	Maize	<i>parviglumis</i>
<i>tb1</i>	<i>mexicana</i>	<i>parviglumis</i>
<i>wip1</i>	<i>mexicana</i>	<i>luxurians</i>
<i>wip1</i>	<i>parviglumis</i>	Maize
<i>wx1</i>	<i>mexicana</i>	Maize

Each line represents a taxon pair/locus combination for which the posterior probability of assignment to the original taxon is <50% for at least one genotype.

Pooling data across loci, our results further suggest the possibility of adaptive evolution in *parviglumis*. The pooled estimate of  $\gamma$  for *parviglumis* is weakly positive and  $NI < 1$  (Table 3); using the  $NI$  values as a crude estimate, as many as 25% of fixed nonsynonymous variants may have been subjected to positive selection. These results are far from conclusive, however, as the estimated value of  $\gamma$  is not significantly different from zero (Table 3) and we do not detect a negative correlation between nonsynonymous divergence and synonymous diversity (Pearson's  $r = -0.25$ ,  $P = 0.39$  after controlling for synonymous divergence), as might be expected from selective sweeps (ANDOLFATTO 2007).

In plant lineages, similar estimates of selection are available only from the genus *Arabidopsis*. Initial investigation found evidence of negative selection in

**TABLE 5**  
**Relative probabilities of alternative divergence models**

Comparison	Model	$K$
<i>parviglumis-luxurians</i>	Isolation	1.00
	Island	0.23
	Sympatry	3.27
<i>parviglumis-mexicana</i>	Allopatry	1.04
	Isolation	1.00
	Island	4.08
<i>parviglumis-maize</i>	Sympatry	1.50
	Allopatry	3.31
	Isolation	1.00
	Island	0.22
	Sympatry	1.17
	Allopatry	0.39

The Bayes factor ( $K$ ) comparing each model to the isolation model is shown.

*Arabidopsis thaliana*, a result that was interpreted as consistent with the species' largely self-pollinating mating system (BUSTAMANTE *et al.* 2002). More recent work, however, has found similar patterns in both *A. thaliana* and its outcrossing relative *A. lyrata* (FOX *et al.* 2008), suggesting that mating system alone is unlikely sufficient to explain differences in estimates of selection. While our data suggest some potential differences in the action of selection among taxa, all four taxa share a similar mating system, and, unlike the predictions of BUSTAMANTE *et al.* (2002), the taxon with the lowest estimated effective population size (*luxurians*) does not exhibit the most negative estimates of  $\gamma$ .

**Models of divergence:** In addition to providing estimates of diversity and recombination, our multi-locus data and stratified sampling of taxa within *Zea* allow for model-based population genetic analysis of divergence. Like the most recent incarnations of the DPG approach (HEY and NIELSEN 2004, 2007; BECQUET and PRZEWORSKI 2007)—and in contrast to other current methodologies (*e.g.*, KUHNER 2006)—our approximate Bayesian approach incorporates both migration and divergence under the isolation model. Unlike other DPG methods, however, our approach allows for both the inclusion of recombination and a framework in which to explicitly test alternative models of population divergence and speciation.

None of our models of divergence can be systematically rejected in all comparisons, and in fact each finds some measure of support ( $K > 1$ ) in at least one of the divergence events compared (Table 5). Overall, our modeling provides strong support for the sympatric model of divergence for the split between *parviglumis* and *luxurians*, for recent gene flow between *mexicana* and *parviglumis* (76% of total probability in the allopatric and island models), and in favor of models devoid of recent introgression for maize and *parviglumis* (78% total probability in the isolation and sympatric models).

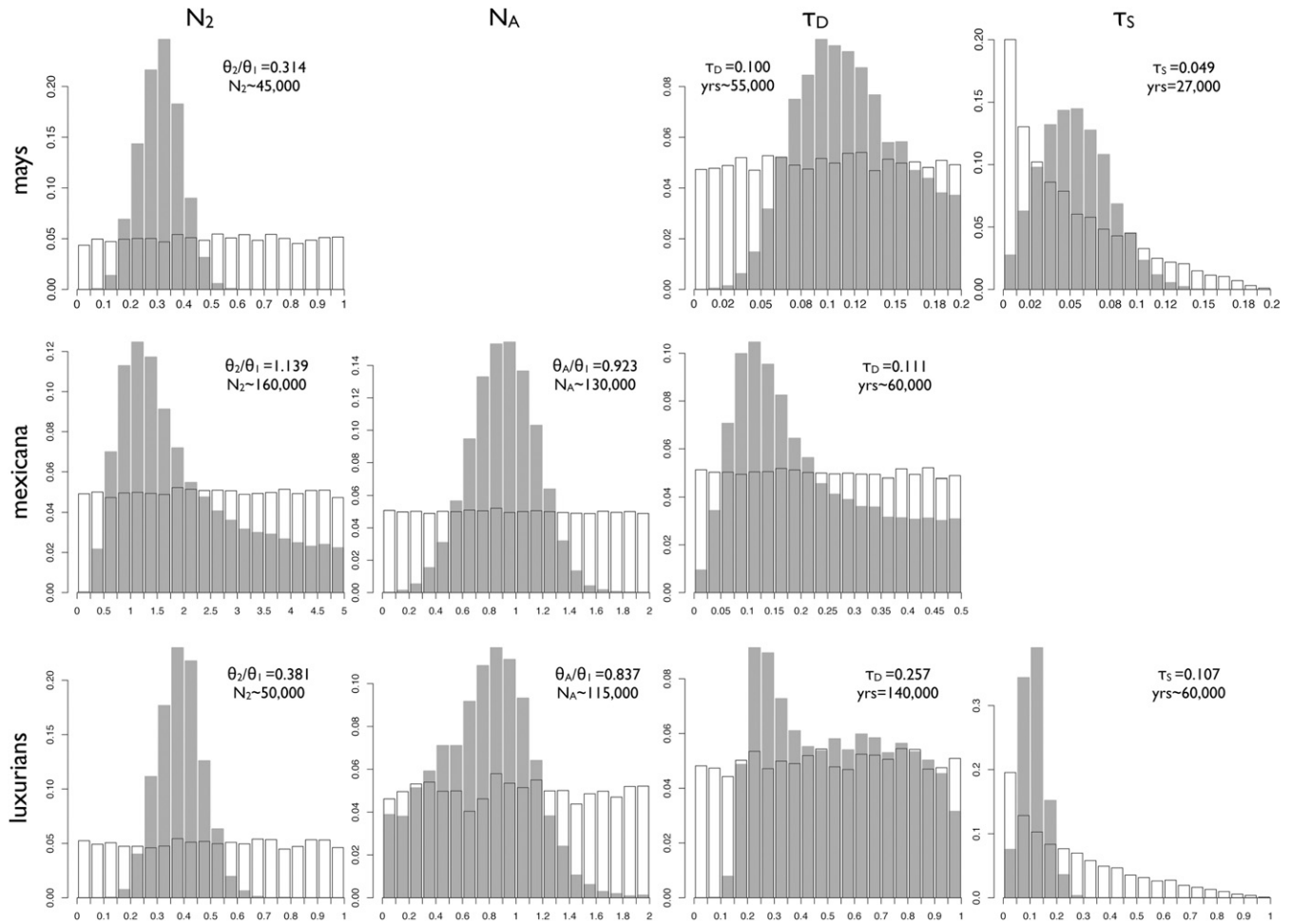


FIGURE 7.—Parameter estimates of the most likely model for each of the three divergence comparisons to *parviglumis*. Posterior probability distributions are shaded and prior probability distributions are open. The mode of each parameter is given along with its conversion into  $N$  or years. Estimates are from the island model for *parviglumis*–*mexicana* and from the sympatric model for *parviglumis*–*luxurians* and *parviglumis*–maize (see Table 5).

Comparison of Bayes factors further hints that gene flow between *mexicana* and *parviglumis* may have been continuous throughout their divergence and that *parviglumis* and maize may have diverged in sympatry. These differences among comparisons suggest that the divergence process is heterogeneous among taxa in *Zea*.

Our Bayesian estimates of divergence time (Figure 7) match surprisingly well with estimates based on allozyme differentiation between *parviglumis* and *luxurians* (135,000 years) and between *parviglumis* and *mexicana* (61,000 years) (HANSON *et al.* 1996). Additionally, our estimates of both divergence time and effective population size are broadly similar to those from previous analyses using fewer loci (HILTON and GAUT 1998; WHITE and DOEBLEY 1999) after accounting for differences in the substitution rates assumed.

Interpretation of the difference between archaeological data (POHL *et al.*, 2007) and our estimate of the *parviglumis*–maize divergence is not as straightforward. One possible explanation is that widespread selection during domestication has caused our loci to appear

more divergent than expected under neutrality. This seems unlikely, however, given the small percentage of the genome thought to be under direct selection during domestication (WRIGHT *et al.* 2005) and the rapid breakdown of linkage disequilibrium in the maize genome (REMINGTON *et al.* 2001). Furthermore, our data lack strong evidence for selection, and reanalysis of the sympatric model excluding loci thought to be important for domestication (*tb1*, *c1*, *d8*, and *ts2*) still provides an overestimate of the divergence time ( $\sim 35,000$  years, data not shown). Alternatively, several other scenarios might explain the high proportion of unique variants seen in maize (Figure 4) and thus the high estimates of divergence time. Our sampling of *parviglumis*, for example, might not have included samples of the population most directly ancestral to domesticated maize, and either a history of population expansion or conscious efforts on the part of cultivators to preserve diversity might have maintained more polymorphism in maize than expected under our constant-size model. Finally, although microsatellite analysis

clearly shows that *parviglumis* is the wild ancestor of maize (MATSUOKA *et al.* 2002), our data highlight the similarities of *mexicana* and maize and the possibility of extensive gene flow between them (see below). Such gene flow may help to explain why the estimated timing of the *parviglumis*–maize divergence is indistinguishable from the timing of the *parviglumis*–*mexicana* divergence with our methods.

**Introgression and its consequences:** Our ABC model comparisons are based on mean patterns of shared, fixed, and exclusive polymorphisms over the entire data set and thus cannot identify gene flow events at individual loci. Additionally, our pairwise modeling approach cannot identify gene flow between taxa that were not compared (*e.g.*, *mexicana*–maize). To complement our modeling efforts, we used STRUCTURE to perform tests for recent gene flow at individual loci. Although a number of loci share identical sequences between one or more of the taxa studied (supplemental Table S7), the STRUCTURE analysis using haplotypes defined by complete sequence identity did not uncover any likely cases of introgression at any locus. This may not be surprising, however; this approach is undoubtedly conservative because the data were reduced to a measure of either complete sequence identity or sequence difference, with no recognition of similarity. Because some individuals shared identity along long regions of sequence, we applied an alternative approach, defining haplotypes within nonrecombining segments (CHEN *et al.* 2009) to represent genotypes at individual loci. With this approach, STRUCTURE analysis identified several sequence segments that appear to suggest recent gene flow (Table 4).

These STRUCTURE results are not without caveats, however. Because each region within a locus was treated as an independent marker, much of the gene flow identified as recent by STRUCTURE could reflect older introgression events broken up by subsequent recombination or may even be the product of shared (but rare) ancestral polymorphism. Moreover, STRUCTURE treats each region as independent, and linkage among these regions could mislead the analyses. Finally, our inferences of historical and recent gene flow—based on ABC and STRUCTURE analyses, respectively—are not in perfect agreement. This is not entirely unexpected, however, because the approaches are fundamentally different. It is not difficult to imagine, for example, that the mean numbers of shared and fixed variants could be suggestive of continued gene flow even if few loci show distinct enough patterns to be identified by STRUCTURE. Similarly, ABC might identify an overall pattern of isolation, but a few individual loci might show patterns suggestive of recent introgression.

In spite of these caveats, inspection of individual alignments supports the plausibility of our inference of recent gene flow for at least some loci. Shared derived indels, identified by comparison to the outgroup

*Tripsacum* but not included in haplotype assignment or clustering analyses, are present in several of the loci highlighted by our STRUCTURE analysis. One particularly striking example is the shared presence of a 131-bp MITE-like transposable element insertion in two alleles of locus *tb1* in *luxurians* and the purportedly introgressed segment in *mexicana*.

Overall, both coalescent and clustering approaches thus provide evidence for gene flow among the taxa of *Zea*. How reasonable are these inferences? While there is little empirical data to support a given model of speciation, previous work investigating introgression and gene flow substantiates our inferences of recent introgression. For example, hybridization is known to occur at least sporadically between domesticated maize and *parviglumis* (WILKES 1977), and, despite genetic incompatibility factors limiting unidirectional gene flow between maize and *mexicana* (EVANS and KERMICLE 2001), the two taxa readily hybridize (BALTAZAR *et al.* 2005; ELLSTRAND *et al.* 2007). Allozyme data further suggest that introgression between maize and *mexicana* may be common (BLANCAS *et al.* 2002), and some evidence suggests that introgression from *mexicana* may have contributed to maize domestication (GALLAVOTTI *et al.* 2004). Finally, although our data support both historical (Table 5) and more recent (Table 4) gene flow between *mexicana* and *parviglumis*, we are not aware of field observations of hybridization between these taxa, and they are not currently found in sympatry. It is likely, however, that the historical distribution of these taxa was different from that observed today (WILKES 1972), and our results are corroborated by microsatellite evidence of admixture between these taxa (FUKUNAGA *et al.* 2005).

Although we identify one possible recent introgression event between *Z. luxurians* and *parviglumis* (Table 4), strong support for the “sympatric model” (Table 5) suggests that gene flow between these taxa has been predominantly historical. This inference is consistent with current geographic patterns, because *luxurians* is isolated from known populations of either *parviglumis* or *mexicana* (WILKES 1977; SÁNCHEZ GONZÁLEZ and RUIZ CORRAL 1997; FUKUNAGA *et al.* 2005). Several lines of evidence nonetheless suggest that our inferred historical introgression is not implausible. Recent collections have identified populations of *luxurians* in the Mexican state of Oaxaca (CUEVAS 2006), not far from the range of *mexicana*. Furthermore, extant populations of a fourth subspecies of *Z. mays*, *ssp. huehuetenangensis*, are currently found in Western Guatemala, suggesting that the ancestral ranges of *Z. mays* and *luxurians* may have overlapped to some extent.

If introgression has played both an historical and a recent role in *Zea*, what are the implications for systematic inference? Clearly, nuclear sequence data from multiple unlinked loci are not ideal for phylogenetic reconstruction because, among other reasons,



introgression can obscure historical relationships. Nonetheless, we can use such data to generate a distance estimate (NEI and LI 1979) that can be converted to divergence time given the rate of substitution. We took this approach to see what light our data could shed on the larger phylogeny of the genus *Zea*. To accomplish this, we made use of sequences of the diploid perennial taxon *Z. diploperennis* (hereafter *diploperennis*) from *adh1*, *c1*, *wip1*, and *wx1* (TIFFIN and GAUT 2001). Nei's net divergence values at these four loci between *parviglumis* and *luxurians* (0.0169) and *parviglumis* and *diploperennis* (0.0148) seem to suggest a more recent common ancestor between *parviglumis* and *diploperennis*. However, using three of these loci to compare *mexicana* to both *luxurians* and *diploperennis* reveals the opposite pattern, and both *luxurians* and *diploperennis* are closer to *Z. mays* than they are to each other (data not shown). These incongruous values echo the conflicting phylogenies found with ribosomal DNA (BUCKLER and HOLTSFORD 1996) and other markers (DOEBLEY *et al.* 1984, 1987; FUKUNAGA *et al.* 2005). The possibility of introgression from *diploperennis* (FUKUNAGA *et al.* 2005) or *luxurians* (DOEBLEY *et al.* 1984; this article) complicates estimates of divergence times from these numbers (TESHIMA and TAJIMA 2002), but they agree well with our Bayesian inference and previous estimates from other markers.

Together, these and other estimates (GAUT and CLEGG 1993; HANSON *et al.* 1996; HILTON and GAUT 1998) suggest that the various species of *Zea* may have arisen nearly contemporaneously—on the order of 100,000–300,000 years using the mutation rate of CLARK *et al.* (2005). Similar divergence measures between our outgroup *Tripsacum* and *parviglumis* (Nei's divergence = 0.0637) or *luxurians* (0.0729) provide estimates of ~1–1.2 million years as the upper limit on the age of the genus. Although this is older than the estimates of WHITE and DOEBLEY (1999) based on divergence at *terminal ear 1* (~0.5 million years using the mutation rate of CLARK *et al.* 2005), it agrees well with the observation from *ITS* data that speciation in *Zea* is recent compared to its divergence from *Tripsacum* (BUCKLER and HOLTSFORD 1996).

Our ABC analyses suggest that these systematic inferences are complicated by historical introgression, and the STRUCTURE analyses open the possibility that recent introgression events may also cloud systematic relationships. One intriguing possibility, which might serve to explain support for recent (but perhaps limited) introgression between all of our taxa (Table 4), is that domesticated maize may serve as a bridge for the flow of alleles among wild taxa. Cultivated maize is nearly ubiquitous in Mexico and Guatemala and can be found in close proximity to all of the wild taxa of *Zea*. Data presented here and elsewhere (WILKES 1977; DOEBLEY *et al.* 1984; BALTAZAR *et al.* 2005; FUKUNAGA *et al.* 2005; ELLSTRAND *et al.* 2007) suggest that maize

exchanges genes with its wild relatives. Broader sampling of maize and its wild relatives will be required to definitively test this hypothesis, but the possibility that domesticated maize may have functioned as a genetic go-between for wild populations certainly has profound implications both for the maintenance *in situ* of teosinte diversity and for our understanding of the risks and consequences of transgene escape from cultivated maize.

We thank J. U'Ren and R. Gaut for sequencing work; K. Thornton for computational assistance; J. J. Sánchez-González for alerting us to the recent collection of *luxurians*; and J. Doebley, K. Thornton, and two anonymous reviewers for helpful discussion. This work was supported by National Science Foundation grant DBI0321467. Part of this work was carried out using the resources of the Computational Biology Service Unit from Cornell University, which is partially funded by Microsoft Corporation

#### LITERATURE CITED

- ANDOLFATTO, P., 2007 Hitchhiking effects of recurrent beneficial amino acid substitutions in the *Drosophila melanogaster* genome. *Genome Res.* **17**: 1755–1762.
- ANTONOVICS, J., 2006 Evolution in closely adjacent plant populations X: long-term persistence of prereproductive isolation at a mine boundary. *Heredity* **97**: 33–37.
- ARNOLD, M. L., 2004 Transfer and origin of adaptations through natural hybridization: Were Anderson and Stebbins right? *Plant Cell* **16**: 562–570.
- ARNOLD, M. L., 2006 *Evolution Through Genetic Exchange*. Oxford University Press, Oxford/New York.
- ARUNYAWAT, U., W. STEPHAN and T. STÄDLER, 2007 Using multilocus sequence data to assess population structure, natural selection, and linkage disequilibrium in wild tomatoes. *Mol. Biol. Evol.* **24**: 2310–2322.
- BALTAZAR, B. M., J. DE JESUS SANCHEZ-GONZALEZ, L. DE LA CRUZ-LARIOS and J. B. SCHOPER, 2005 Pollination between maize and teosinte: an important determinant of gene flow in Mexico. *Theor. Appl. Genet.* **110**: 519–526.
- BECQUET, C., and M. PRZEWORSKI, 2007 A new approach to estimate parameters of speciation models with application to apes. *Genome Res.* **17**: 1505–1519.
- BLANCAS, L., D. I. ARIAS and N. C. ELLSTRAND, 2002 Patterns of genetic diversity in sympatric and allopatric populations of maize and its wild relative teosinte in Mexico: evidence for hybridization, pp. 31–38 in *Scientific Methods Workshop: Ecological and Agronomic Consequences of Gene Flow From Transgenic Crops to Wild Relatives*, edited by A. SNOW. Ohio State University, Columbus, OH.
- BUCKLER, E. S. T., and T. P. HOLTSFORD, 1996 *Zea* systematics: ribosomal ITS evidence. *Mol. Biol. Evol.* **13**: 612–622.
- BUSTAMANTE, C. D., R. NIELSEN, S. A. SAWYER, K. M. OLSEN, M. D. PURUGGANAN *et al.*, 2002 The cost of inbreeding in *Arabidopsis*. *Nature* **416**: 531–534.
- CHEN, H., P. L. MORRELL, V. E. ASHWORTH, M. DE LA CRUZ and M. T. CLEGG, 2009 Tracing the geographic origins of major avocado cultivars. *J. Hered.* **100**: 56–65.
- CLARK, A. G., 1997 Neutral behavior of shared polymorphism. *Proc. Natl. Acad. Sci. USA* **94**: 7730–7734.
- CLARK, R. M., S. TAVARE and J. DOEBLEY, 2005 Estimating a nucleotide substitution rate for maize from polymorphism at a major domestication locus. *Mol. Biol. Evol.* **22**: 2304–2312.
- COYNE, J. A., and H. A. ORR, 2004 *Speciation*. Sinauer Associates, Sunderland, MA.
- CUEVAS, F. A., 2006 Nueva población de teocintle en Oaxaca. XXI Congreso Nacional y Primero Internacional de Fitogenética, Sociedad Mexicana de Fitogenética, Tuxtla Gutiérrez, México.
- DOEBLEY, J. F., 1990 Molecular evidence for gene flow among *Zea* species. *Bioscience* **40**: 443–448.

- DOEBLEY, J. F., 2004 The genetics of maize evolution. *Annu. Rev. Genet.* **38**: 37–59.
- DOEBLEY, J. F., and H. H. ILTIS, 1980 Taxonomy of *Zea* (Gramineae). 1. A subgeneric classification with key to taxa. *Am. J. Bot.* **67**: 982–993.
- DOEBLEY, J. F., M. M. GOODMAN and C. W. STUBER, 1984 Isoenzymatic variation in *Zea* (Gramineae). *Syst. Bot.* **9**: 203–218.
- DOEBLEY, J. F., W. RENFROE and A. BLANTON, 1987 Restriction site variation in the *Zea* chloroplast genome. *Genetics* **117**: 139–147.
- ELLSTRAND, N. C., L. C. GARNER, S. HEGDE, R. GUADAGNUOLO and L. BLANCAS, 2007 Spontaneous hybridization between maize and teosinte. *J. Hered.* **98**: 183–187.
- EVANS, M. M. S., and J. L. KERMICLÉ, 2001 Teosinte crossing barrier 1, a locus governing hybridization of teosinte with maize. *Theor. Appl. Genet.* **103**: 259–265.
- EYRE-WALKER, A., R. L. GAUT, H. HILTON, D. L. FELDMAN and B. S. GAUT, 1998 Investigation of the bottleneck leading to the domestication of maize. *Proc. Natl. Acad. Sci. USA* **95**: 4441–4446.
- FAGUNDES, N. J. R., N. RAY, M. BEAUMONT, S. NEUENSCHWANDER, F. M. SALZANO *et al.*, 2007 Statistical evaluation of alternative models of human evolution. *Proc. Natl. Acad. Sci. USA* **104**: 17614–17619.
- FOXÉ, J. P., V. U. DAR, H. ZHENG, M. NORDBORG, B. S. GAUT *et al.*, 2008 Selection on amino acid substitutions in *Arabidopsis*. *Mol. Biol. Evol.* **25**: 1375–1383.
- FU, Y. X., and W. H. LI, 1993 Statistical test of neutrality of mutations. *Genetics* **133**: 693–709.
- FUKUNAGA, K., J. HILL, Y. VIGOUROUX, Y. MATSUOKA, G. J. SANCHEZ *et al.*, 2005 Genetic diversity and population structure of teosinte. *Genetics* **169**: 2241–2254.
- GALLAVOTTI, A., Q. ZHAO, J. KYOZUKA, R. B. MEELEY, M. K. RITTER *et al.*, 2004 The role of barren stalk1 in the architecture of maize. *Nature* **432**: 630–635.
- GAUT, B. S., and M. T. CLEGG, 1993 Molecular evolution of the *Adh1* locus in the genus *Zea*. *Proc. Natl. Acad. Sci. USA* **90**: 5095–5099.
- GAUT, B. S., B. R. MORTON, B. C. MCCAIG and M. T. CLEGG, 1996 Substitution rate comparisons between grasses and palms: synonymous rate differences at the nuclear gene *Adh* parallel rate differences at the plastid gene *rbcl*. *Proc. Natl. Acad. Sci. USA* **93**: 10274–10279.
- HAMRICK, J. L., and M. J. W. GODT, 1989 Allozyme diversity in plant species, pp. 43–63 in *Plant Population Genetics, Breeding, and Genetic Resources*, edited by A. H. D. BROWN, M. T. CLEGG, A. L. KAHLER and B. S. WEIR. Sinauer, Sunderland, MA.
- HAMRICK, J. L., and J. D. NASON, 1996 Consequences of dispersal in plants, pp. 203–236 in *Population Dynamics in Ecological Space and Time*, edited by O. E. RHODES, R. K. CHESSER and M. H. SMITH. University of Chicago Press, Chicago.
- HANSON, M. A., B. S. GAUT, A. O. STEC, S. I. FURSTENBERG, M. M. GOODMAN *et al.*, 1996 Evolution of anthocyanin biosynthesis in maize kernels: the role of regulatory and enzymatic loci. *Genetics* **143**: 1395–1407.
- HEY, J., and R. NIELSEN, 2004 Multilocus methods for estimating population sizes, migration rates and divergence time, with applications to the divergence of *Drosophila pseudoobscura* and *D. persimilis*. *Genetics* **167**: 747–760.
- HEY, J., and R. NIELSEN, 2007 Integration within the Felsenstein equation for improved Markov chain Monte Carlo methods in population genetics. *Proc. Natl. Acad. Sci. USA* **104**: 2785–2790.
- HILTON, H., and B. S. GAUT, 1998 Speciation and domestication in maize and its wild relatives: evidence from the globulin-1 gene. *Genetics* **150**: 863–872.
- ILTIS, H. H., and J. F. DOEBLEY, 1980 Taxonomy of *Zea* (Gramineae). 2. Subspecific categories in the *Zea mays* complex and a generic synopsis. *Am. J. Bot.* **67**: 994–1004.
- HUDSON, R. R., 1990 Gene genealogies and the coalescent process. *Oxf. Surv. Evol. Biol.* **7**: 1–44.
- HUDSON, R. R., and N. L. KAPLAN, 1985 Statistical properties of the number of recombination events in the history of a sample of DNA sequences. *Genetics* **111**: 147–164.
- JEFFREYS, H., 1998 *Theory of Probability*. Oxford University Press, Oxford.
- KUHNER, M. K., 2006 LAMARC 2.0: maximum likelihood and Bayesian estimation of population parameters. *Bioinformatics* **22**: 768–770.
- MALLET, J., 2007 Hybrid speciation. *Nature* **446**: 279–283.
- MARJORAM, P., and S. TAVARE, 2006 Modern computational approaches for analyzing molecular genetic variation data. *Nat. Rev. Genet.* **7**: 759–770.
- MATSUOKA, Y., Y. VIGOUROUX, M. M. GOODMAN, G. J. SANCHEZ, E. BUCKLER *et al.*, 2002 A single domestication for maize shown by multilocus microsatellite genotyping. *Proc. Natl. Acad. Sci. USA* **99**: 6080–6084.
- MCDONALD, J. H., and M. KREITMAN, 1991 Adaptive protein evolution at the *Adh* locus in *Drosophila*. *Nature* **351**: 652–654.
- MCVEAN, G., P. AWADALLA and P. FEARNEHEAD, 2002 A coalescent-based method for detecting and estimating recombination from gene sequences. *Genetics* **160**: 1231–1241.
- MOELLER, D. A., M. I. TENAILLON and P. TIFFIN, 2007 Population structure and its effects on patterns of nucleotide polymorphism in teosinte (*Zea mays* ssp. *parviglumis*). *Genetics* **176**: 1799–1809.
- MORJAN, C. L., and L. H. RIESEBERG, 2004 How species evolve collectively: implications of gene flow and selection for the spread of advantageous alleles. *Mol. Ecol.* **13**: 1341–1356.
- NEI, M., and W. H. LI, 1979 Mathematical model for studying genetic variation in terms of restriction endonucleases. *Proc. Natl. Acad. Sci. USA* **76**: 5269–5273.
- POHL, M. E., D. R. PIPERNO, K. O. POPE and J. G. JONES, 2007 Microfossil evidence for pre-Columbian maize dispersals in the neotropics from San Andres, Tabasco, Mexico. *Proc. Natl. Acad. Sci. USA* **104**: 6870–6875.
- PRITCHARD, J. K., M. STEPHENS and P. DONNELLY, 2000 Inference of population structure using multilocus genotype data. *Genetics* **155**: 945–959.
- RAMOS-ONSINS, S. E., B. E. STRANGER, T. MITCHELL-OLDS and M. AGUADE, 2004 Multilocus analysis of variation and speciation in the closely related species *Arabidopsis halleri* and *A. lyrata*. *Genetics* **166**: 373–388.
- RAND, D. M., and L. M. KANN, 1996 Excess amino acid polymorphism in mitochondrial DNA: contrasts among genes from *Drosophila*, mice, and humans. *Mol. Biol. Evol.* **13**: 735–748.
- REMINGTON, D. L., J. M. THORNSBERRY, Y. MATSUOKA, L. M. WILSON, S. R. WHITT *et al.*, 2001 Structure of linkage disequilibrium and phenotypic associations in the maize genome. *Proc. Natl. Acad. Sci. USA* **98**: 11479–11484.
- RIESEBERG, L. H., and J. H. WILLIS, 2007 Plant speciation. *Science* **317**: 910–914.
- SÁNCHEZ GONZÁLEZ, J. J., and J. A. RUIZ CORRAL, 1997 Teosinte distribution in Mexico, pp. 18–36 in *Gene Flow Among Maize Landraces, Improved Maize Varieties and Teosinte: Implications for Transgenic Maize*, edited by J. A. SERRATOS, M. C. WILLCOX and F. CASTILLO GONZÁLEZ. CIMMYT, Mexico City.
- SAVOLAINEN, V., M. C. ANSTETT, C. LEXER, I. HUTTON, J. J. CLARKSON *et al.*, 2006 Sympatric speciation in palms on an oceanic island. *Nature* **441**: 210–213.
- STÄDLER, T., K. ROSELIUS and W. STEPHAN, 2005 Genealogical footprints of speciation processes in wild tomatoes: demography and evidence for historical gene flow. *Evolution* **59**: 1268–1279.
- STÄDLER, T., U. ARUNYAWAT and W. STEPHAN, 2008 Population genetics of speciation in two closely related wild tomatoes (*Solanum* section *Lycopersicon*). *Genetics* **178**: 339–350.
- STEBBINS, G. L., 1981 Coevolution of grasses and herbivores. *Ann. Mo. Bot. Gard.* **68**: 75–86.
- TAJIMA, F., 1989 Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics* **123**: 585–595.
- TENAILLON, M. I., M. C. SAWKINS, A. D. LONG, R. L. GAUT, J. F. DOEBLEY *et al.*, 2001 Patterns of DNA sequence polymorphism along chromosome 1 of maize (*Zea mays* ssp. *mays* L.). *Proc. Natl. Acad. Sci. USA* **98**: 9161–9166.
- TENAILLON, M. I., J. U'REN, O. TENAILLON and B. S. GAUT, 2004 Selection versus demography: a multilocus investigation of the domestication process in maize. *Mol. Biol. Evol.* **21**: 1214–1225.
- TESHIMA, K. M., and F. TAJIMA, 2002 The effect of migration during the divergence. *Theor. Popul. Biol.* **62**: 81–95.
- THORNTON, K., 2003 libsequence: a C++ class library for evolutionary genetic analysis. *Bioinformatics* **19**: 2325–2327.
- TIFFIN, P., and B. S. GAUT, 2001 Sequence diversity in the tetraploid *Zea perennis* and the closely related diploid *Z. diploperennis*: insights from four nuclear loci. *Genetics* **158**: 401–412.

- VIGOUROUX, Y., S. MITCHELL, Y. MATSUOKA, M. HAMBLIN, S. KRESOVICH *et al.*, 2005 An analysis of genetic diversity across the maize genome using microsatellites. *Genetics* **169**: 1617–1630.
- WAKELEY, J., and J. HEY, 1997 Estimating ancestral population parameters. *Genetics* **145**: 847–855.
- WANG, R. L., J. WAKELEY and J. HEY, 1997 Gene flow and natural selection in the origin of *Drosophila pseudoobscura* and close relatives. *Genetics* **147**: 1091–1106.
- WATTERSON, G. A., 1975 On the number of segregating sites in genetical models without recombination. *Theor. Popul. Biol.* **7**: 256–276.
- WHITE, S. E., and J. F. DOEBLEY, 1999 The molecular evolution of terminal ear 1, a regulatory gene in the genus Zea. *Genetics* **153**: 1455–1462.
- WILKES, H. G., 1972 Maize and its wild relatives. *Science* **177**: 1071–1077.
- WILKES, H. G., 1977 Hybridization of maize and teosinte, in Mexico and Guatemala and the improvement of maize. *Econ. Bot.* **31**: 254–293.
- WOLFE, K. H., W. H. LI and P. M. SHARP, 1987 Rates of nucleotide substitution vary greatly among plant mitochondrial, chloroplast, and nuclear DNAs. *Proc. Natl. Acad. Sci. USA* **84**: 9054–9058.
- WRIGHT, S. I., I. V. BI, S. G. SCHROEDER, M. YAMASAKI, J. F. DOEBLEY *et al.*, 2005 The effects of artificial selection on the maize genome. *Science* **308**: 1310–1314.
- ZHANG, L. B., and S. GE, 2007 Multilocus analysis of nucleotide variation and speciation in *Oryza officinalis* and its close relatives. *Mol. Biol. Evol.* **24**: 769–783.

Communicating editor: D. CHARLESWORTH