

## Supplemental Methods

### *Modeling expected crossovers*

To model expected crossovers, an F1 chromosome with 10,000 equally spaced sites was constructed by coding one parental gamete as a sequence of 0's and the other as a sequence of 1's. To simulate a selfing generation, two gametes were generated by randomly generating the position of a single crossover from a uniform distribution. Because only one chromatid of each chromosome of a pair participates in a crossover, there is a .5 probability that a generated gamete will contain a crossover. Two independently generated gametes were combined to create a self-progeny. The selfing process was repeated to simulate multiple generations of selfing. 100,000 progeny were simulated and the resulting number and type of crossovers were counted. Because each generated gamete has a 0.5 probability of containing a cross-over, it represents a chromosome segment with a genetic map length of 50cM. The resulting crossover rate was multiplied by 28 to calculate the rate expected for a 1400 cM map length.

### *Quality control of imputed genotypes*

After an initial imputation of RIL genotypes, RILs that had more than 30% heterozygous loci were excluded from the analysis. A likely reason for apparent excess heterozygosity was the presence of non-parental haplotypes resulting from contaminated pollinations or sample mixups during the inbreeding process. After examining several alternative quality parameters, it was found that excess heterozygosity most consistently identified problem genotypes. After filtering RILs on genotype quality, 4714 US-NAM RILs and 1382 CN-NAM RILs remained and were re-imputed. Following that, intervals between adjacent homozygous sites from different parents were identified as homozygous crossovers. Intervals between heterozygous sites and homozygous sites were identified as heterozygous crossovers. Subsequent analysis used only the homozygous cross-overs (**See Supplemental Results**). Results from the GBS pipeline were validated by comparing them to low density co-dominant markers that were previously used to genotype the US-NAM population (1). We observed that homozygous segments closely matched in the two data sets, but noticed that heterozygous segments in the imputed GBS data sometimes had short homozygous segments embedded in them, resulting in too many heterozygous crossovers called. Rather than work on tuning the Viterbi algorithm to reduce that, we chose to eliminate the problem by only working with homozygous crossovers. Sampling variability resulting from bulking multiple plants to create DNA samples and contamination of a low frequency of samples with non-parental DNA from stray pollen made consistent imputation of heterozygous segments challenging. Fortunately, these effects showed up as excess heterozygous crossovers and could therefore be largely eliminated.

### *Estimating hotspot significance*

We assessed the significance of the identified hotspots by determining the number of hotspots expected under the observed large-scale (1-Mb) pattern of crossovers if there were no recombination hotspots present (i.e. the probability of crossover were uniform across a 1-Mb interval). In order to generate the null distribution, we divided each chromosome into 1Mb windows, calculated the number of crossovers that occurred in each window within US-NAM, and set the probability of a crossover occurring within a given window proportional to the empirical count. Then, for each RIL, we sampled, with replacement, a number of windows equal to the observed number of homozygous crossover for that RIL. The crossover positions were then sampled from a uniform distribution within the sampled windows, and the nearest flanking GBS markers for the sampled position were used to define the crossover intervals.

### *Whole genome alignment of monocots, GERP score, and phastBias calculation*

For the whole genome alignment, we used the LASTZ/MULTIZ approach adopted by the UCSC genome browser and previously used to align 20 angiosperm genomes to the *A. thaliana* reference (2). Briefly, we aligned the Repeatmasked genomes of *Sorghum bicolor*, *Oryza sativa*, *Setaria italica*, *Hordeum vulgare*, *Brachypodium distachyon*, and *Musa acuminata* to the *Z. mays* B73 version 2 genome. Each query genome was split into 1Mb segments and aligned to each of the B73 chromosomes. We then joined these pairwise alignments into larger chains using axtChain (3), found the best chained alignment using chainNet, and converted these to multiple alignment files. The individual pairwise alignments were then joined into a multiple sequence alignment using the roast program in the MULTIZ package(4).

Following the alignment, we calculated GERP scores for each site in the *Z. mays* genome using the GERP++ package (5). We used 4-fold degenerate sites in the *Z. mays* genome to create the neutral reference tree, which was calculated with approximate maximum likelihood using FastTree2 (6). In order to maintain comparable GERP rates across all regions and limit ourselves to the most robust estimates, we limited further analyses involving GERP to sites at which all species were able to be aligned to the *Z. mays* reference. We also calculated the posterior probability of biased gene conversion at each site using the phastBias program (7), part of the PHAST package.

### *Identification and analysis of hotspot-enriched sequence motifs*

We randomly split hotspots with at least 10 comparison controls into training (n=377) and testing (n=111) sets. Within each set, the hotspots provided the positive examples, while a control region was sampled from each hotspot to form the negative examples. Motifs significantly enriched within the positive training set relative to the negative training set were then identified using Discriminative DNA Motif Discovery (DREME), part of the MEME suite (8). Significantly enriched motifs in the training set were then tested for enrichment within the testing set using Fisher's Exact Test.

Following identification of enriched motifs, we tested their predictive power with respect to recombination in 30kb windows. We created a set of 20,000 of these 30kb windows, which did not overlap with each other or any of the regions used for motif discovery. Of these 20,000 regions, 15,000 were randomly chosen as a training set with 5,000 reserved for testing. We then used stochastic gradient boosting with absolute error loss to predict recombination frequency using the motifs with p-values  $\leq 0.05$  in the discovery test set with and without methylation included. Stochastic gradient boosting was chosen for its generally high performance as an "out-of-the-box" learner (9). The performance of the classifiers was assessed using the root-mean-square error (RMSE), while the relative importance of predictor variables was determined by permuting each variable 1,000 times and calculating the mean reduction in RMSE.

## **Supplemental Results**

### *Comparison to prior genotyping results*

We compared the results of this study to an earlier analysis of the US-NAM population using 1106 markers scored with an Illumina GoldenGate (GG) genotyping assay (1). Since that time some additional lines were added to the population to replace problem lines, and the seed stocks have been increased by sib-pollination. The first study reported that about 136,000 crossovers were detected in 4699 RILs. A reanalysis of the data showed that these consisted of 101,022 homozygous recombinants and

34,666 heterozygous recombinants. The study reported here detected 103,459 homozygous recombinants and 92,276 heterozygous recombinants in 4714 RILs. The rate of homozygous crossovers per RIL was 21.50 and 21.95 for the earlier and current studies, respectively. The rate of heterozygous crossovers was 7.37 and 19.7.

Correctly scoring heterozygous loci was challenging using the GG assay because the heterozygotes did not form a distinct cluster for all markers. To deal with that problem, ambiguous calls were set to missing. Scoring heterozygous loci with GBS is hindered by low read depth, where often only one of the alleles present in a sample will be observed. The GBS problem was handled by using an HMM to call genotypes. To complicate matters, the DNA samples were created by bulking tissue from 4 plants derived from the original S5 plant with one or more generations of self or sib pollination between the original plant and the sampled plants. Because of sampling variation, the ratio of alleles in a sample at a site that is heterozygous in the S5 plant can vary substantially from the expected 1:1 ratio, making heterozygous loci more difficult to score correctly by either method. An additional challenge can result from low levels of foreign DNA contamination, which may be imputed as heterozygous loci. The difficulty of accurately scoring heterozygotes may be partially responsible for the excess of heterozygous crossovers in the imputed GBS data, where spurious homozygous calls in heterozygous segments may appear as multiple crossovers. Conversely, undercalling of heterozygous crossovers in the GG data due to setting heterozygous segments to missing can give rise to the same effect.

In order to compare crossover locations derived from the GBS data to locations derived from the GG data, we used physical positions for the GG markers from the Panzea database (<http://www.panzea.org>), which were determined by BLASTing the sequences used to design the microarray against version 2 of the B73 reference genome. Approximately 95% of the GG SNPs had direct blast hits. The positions of the remainder were either estimated based on the genetic map position or using nearby SNPs from the same amplicon. For each RIL, a crossover interval from the GBS data was taken to be congruent to a GG-based interval if the intervals overlapped.

Comparisons of the crossover intervals of GBS and GG show that approximately 99% of homozygous GG intervals are also identified using GBS, though approximately 15% of these were classified as heterozygous using the GBS method. 83% of the GG heterozygous intervals contained GBS, which were mostly heterozygous. Conversely, 85% of GBS homozygous crossovers fell within GG crossover intervals and 62% of GBS heterozygous crossovers fell within GG crossover intervals.

The concordance of GBS and GG crossovers clearly demonstrates that heterozygous crossovers are not imputed as reliably as homozygous crossovers in either the GBS or the GG data. Furthermore, consideration of the HMM-based imputation method shows that the position of heterozygous crossovers is often less precise. For example, if the parental alleles are coded A and B, and the sequence ABABABABAB is observed, then the sequence is interpreted by the Viterbi algorithm as heterozygous with the observation of a single random parental allele at each site. As such, the sequence ABABABABAAAAAAA would be inferred to contain a heterozygous crossover. However, there is uncertainty governing whether the first A in the A-series is homozygous or heterozygous. The second A is more likely to come from a homozygous locus but could still be a randomly sampled allele from a heterozygous locus.

Given the lower level of confidence for the heterozygous crossover calls and the relative imprecision of their positions, we elected not to include them in further analyses. However, we find that their removal is unlikely to bias our results. The densities of

homozygous and heterozygous crossovers are nearly identical across the entire genome (**Supplemental Figure S1**). Moreover, the estimation of crossover enrichment on the narrow-scale is not significantly altered by the inclusion of heterozygous crossovers in the dataset (Pearson correlation = 0.964) (**Supplemental Figure S34**). Therefore, while we believe that the removal of heterozygous crossovers is the appropriate conservative approach, our results are robust to their inclusion.

### Supplemental References

1. McMullen MD, et al (2009) Genetic properties of the maize nested association mapping population. *Science* 325(5941): 737-740.
2. Hupaló D & Kern AD (2013) Conservation and functional element discovery in 20 angiosperm plant genomes. *Mol Biol Evol* 30(7): 1729-1744.
3. Kent WJ, Baertsch R, Hinrichs A, Miller W & Haussler D (2003) Evolution's cauldron: Duplication, deletion, and rearrangement in the mouse and human genomes. *Proc Natl Acad Sci U S A* 100(20): 11484-11489.
4. Blanchette M, et al (2004) Aligning multiple genomic sequences with the threaded blockset aligner. *Genome Res* 14(4): 708-715.
5. Davydov EV, et al (2010) Identifying a high fraction of the human genome to be under selective constraint using GERP++. *PLoS Comput Biol* 6(12): e1001025.
6. Price MN, Dehal PS & Arkin AP (2010) FastTree 2--approximately maximum-likelihood trees for large alignments. *PLoS One* 5(3): e9490.
7. Capra JA, Hubisz MJ, Kostka D, Pollard KS & Siepel A (2013) A model-based analysis of GC-biased gene conversion in the human and chimpanzee genomes. *PLoS Genet* 9(8): e1003684.
8. Bailey TL (2011) DREME: motif discovery in transcription factor ChIP-seq data. *Bioinformatics* 27(12): 1653-1659.
9. Friedman, JH (2002) Stochastic gradient boosting. *Computational Statistics & Data Analysis* 38(4): 367-378.

*Supplemental Tables*

**Table S1:** Coefficients for terms in linear model of homozygous crossover density in 1Mb windows without inclusion of GBS density. Note that all explanatory variables have been centered and scaled to have a standard deviation of 1.

	Estimate	SE	SS	F-value	p-value
Telomere	-0.0976	0.0104	12.191	87.885	$< 2 \times 10^{-16}$
CpG	-0.4502	0.0226	55.0	396.727	$< 2 \times 10^{-16}$
CHH	0.1096	0.0182	5.0	36.278	$2.023 \times 10^{-9}$
CpG:CHH	0.0872	0.0073	31.64	141.387	$< 2 \times 10^{-16}$
GC	0.0616	0.0100	5.291	38.145	$7.90 \times 10^{-10}$
Repeat	-0.2604	0.0180	29.006	209.104	$< 2 \times 10^{-16}$

Cross-Validation R2: 0.8225

**Table S2:** Coefficients for terms in linear model for CpG methylation in hotspots

	Estimate	SE	SS	F-value	p-value
GBS	-22.690	0.03933	0.18688	27.739	$2.44 \times 10^{-7}$
Crossover enrichment	-0.0451	0.00868	0.18205	27.022	$3.44 \times 10^{-7}$
GBS:Crossover enrichment	3.272	0.86822	0.09571	14.206	0.000193

Adjusted R<sup>2</sup>: 0.231

**Table S3:** Coefficients for terms in linear model for CHG methylation in hotspots

	Estimate	SE	SS	F-value	p-value
GBS	-12.869	3.211	0.0611	16.067	$7.51 \times 10^{-5}$
Crossover enrichment	-0.0204	0.00647	0.03715	9.9262	0.00177
GBS:Crossover enrichment	1.798	0.64713	0.02888	7.162	0.00577

Adjusted R<sup>2</sup>: 0.1363

**Table S4:** Enrichment within a testing set of hotspots for motifs found to be significantly enriched within a distinct training set of hotspots

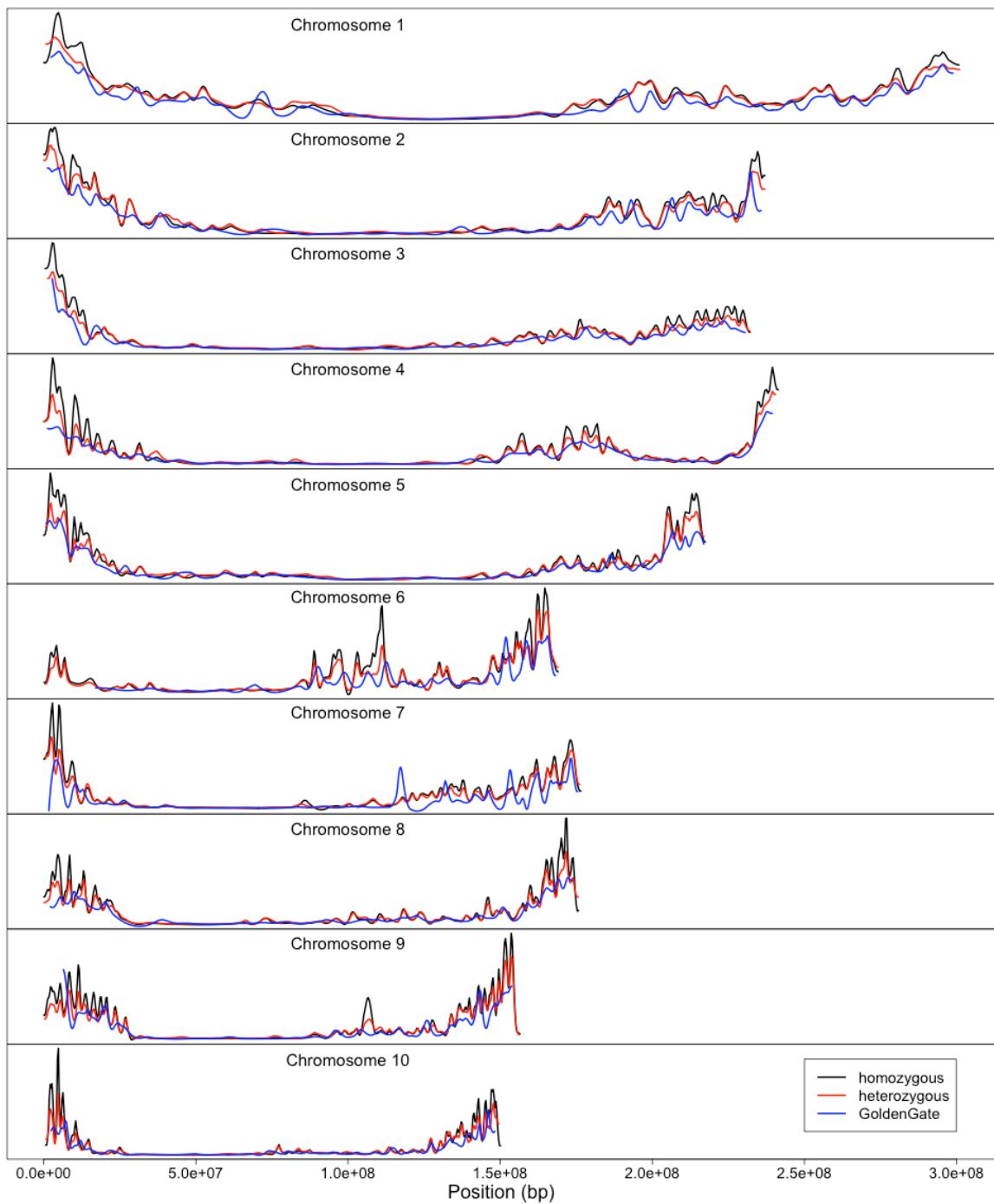
<b>Motif</b>	<b>Hotspot Frequency</b>	<b>Control Frequency</b>	<b>p-value</b>
CGTGACR	76/111	50/111	0.000670
MCGATCGA	43/111	26/111	0.009992
CGTACGTR	70/111	52/111	0.010808
GTGCGTGS	67/111	50/111	0.015628
CCGGCCSS	97/111	84/111	0.018531
CACGCACK	62/111	46/111	0.021873
ACGGSGGC	86/111	72/111	0.026811
CGCGCGCS	85/111	71/111	0.027919
ACGMGACG	63/111	48/111	0.029998
CGCGTBGC	81/111	68/111	0.043057
CCGGCGCH	87/111	75/111	0.047990
CGTASTAC	43/111	33/111	0.101433
GCTAGCTA	67/111	59/111	0.171497
GCTAGKAC	36/111	29/111	0.188132
ACGACGGY	72/111	65/111	0.203750
ACGTACWG	31/111	25/111	0.219944
GGCASGCA	76/111	70/111	0.239774
ATSGATCG	34/111	31/111	0.384076
ACGCTRCG	37/111	35/111	0.443032
CGCSAGCW	82/111	83/111	0.620590
CGTGCKC	41/111	45/111	0.754487

**Table S5:** Coefficients for terms in linear model for GC content in hotspots

	<b>Estimate</b>	<b>SE</b>	<b>SS</b>	<b>F-value</b>	<b>p-value</b>
GBS	3.449	0.2966	0.15417	135.164	$< 2 \times 10^{-16}$
Crossover enrichment	0.0115	0.0017	0.05205	45.633	$3.64 \times 10^{-11}$

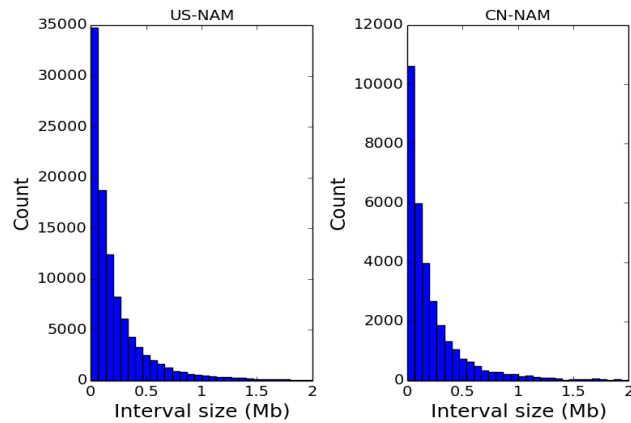
Adjusted  $R^2$ : 0.2863

Supplemental Figures

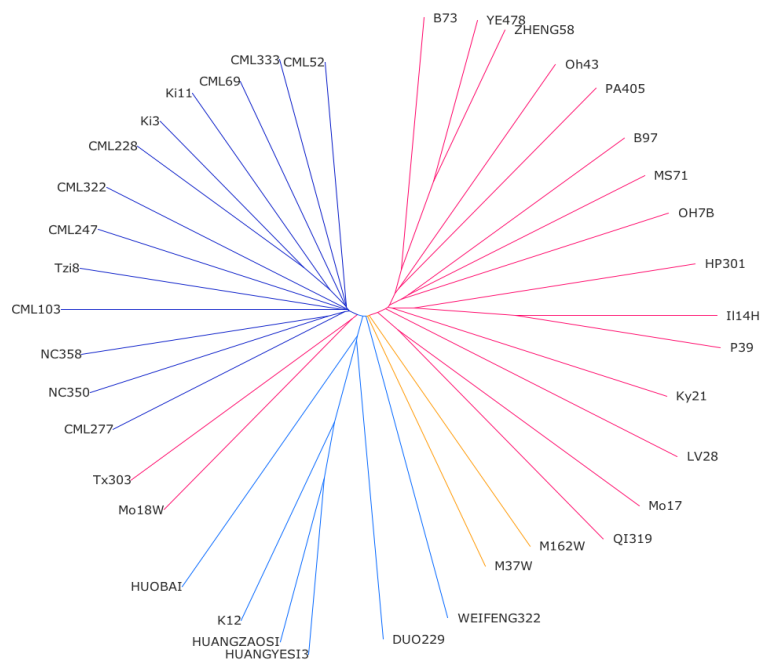


**Figure S1.** Crossover density of homozygous and heterozygous crossovers in the current study, along with the total crossover density in a previous GoldenGate assay of NAM recombination

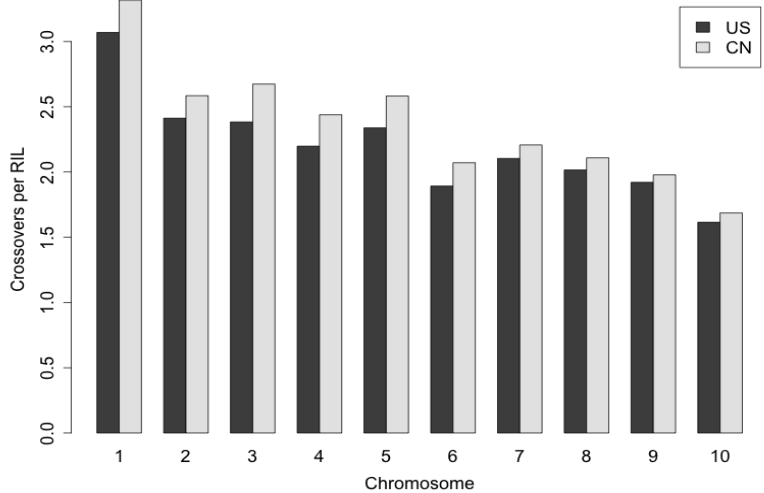




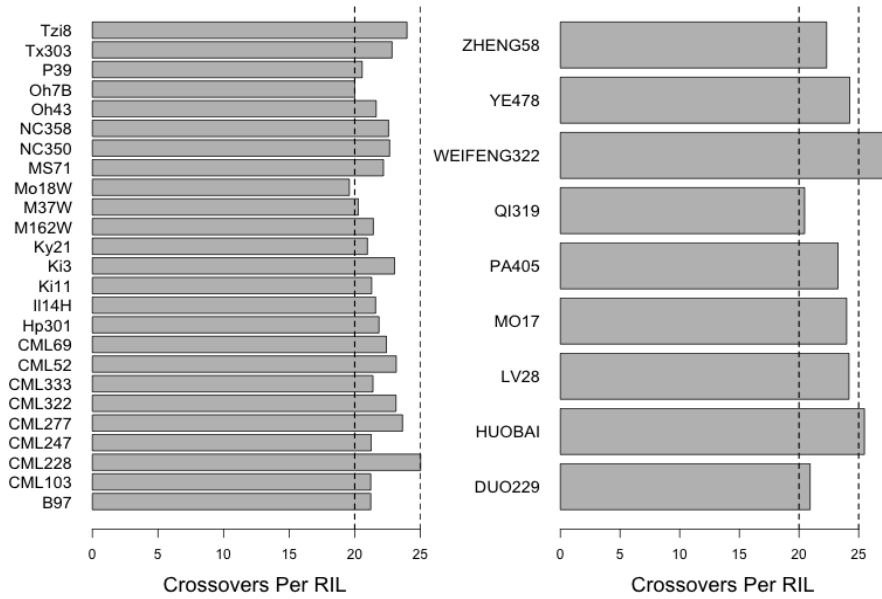
**Figure S2.** Sizes of crossover intervals in US and CN NAM populations



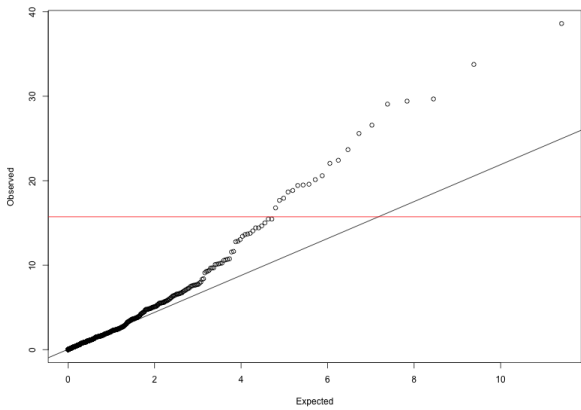
**Figure S3.** Neighbor joining tree for US-NAM and CN-NAM founders. Branch colors are red: founders developed in the US or CN founders that cluster with US lines; purple: tropical founders; blue: CN founders not clustering with US lines; yellow: lines developed in South Africa. Tx303 and Mo18W were developed in the US from tropical parents and cluster with the tropical founders.



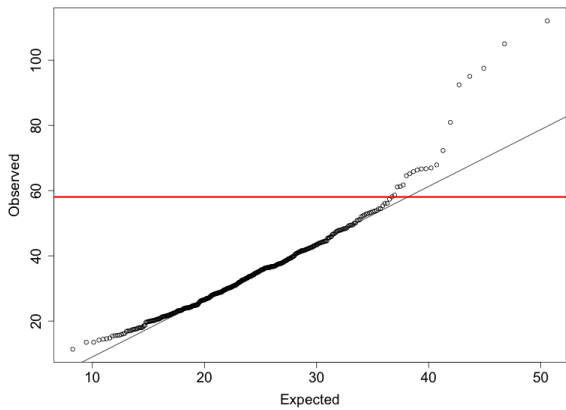
**Figure S4.** Homozygous crossover counts for US-NAM and CN-NAM by chromosome



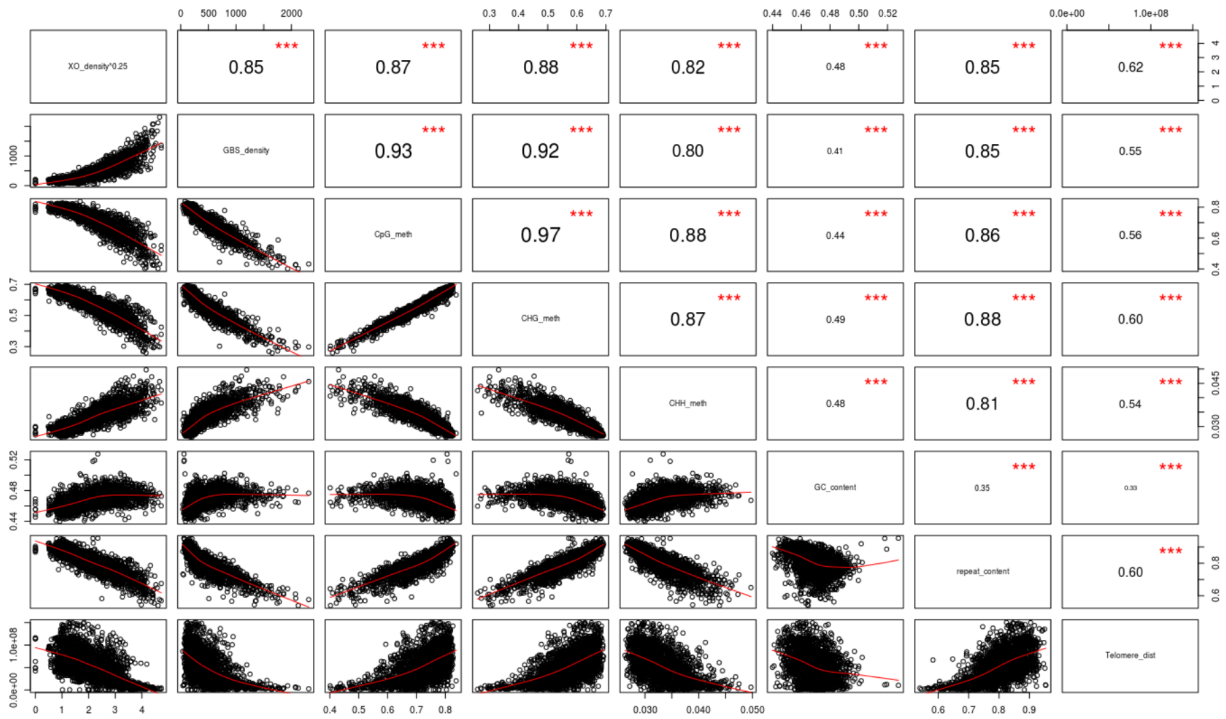
**Figure S5.** Homozygous crossovers per RIL by family for US-NAM (left) and CN-NAM (right). Dashed lines show that most counts fall between 20 and 25.



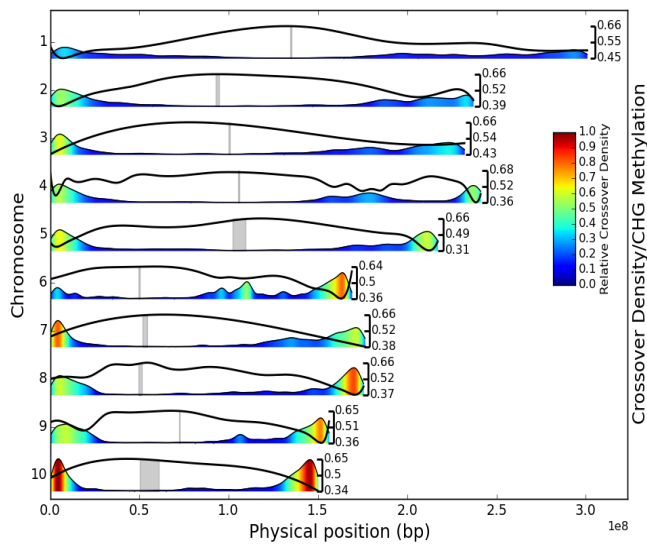
**Figure S6.** QQ-plot of  $\chi^2$  statistic testing equal numbers of crossovers in US-NAM vs. CN-NAM. The red line gives the threshold defined by the Bonferroni-corrected type I error rate,  $\alpha=0.05$



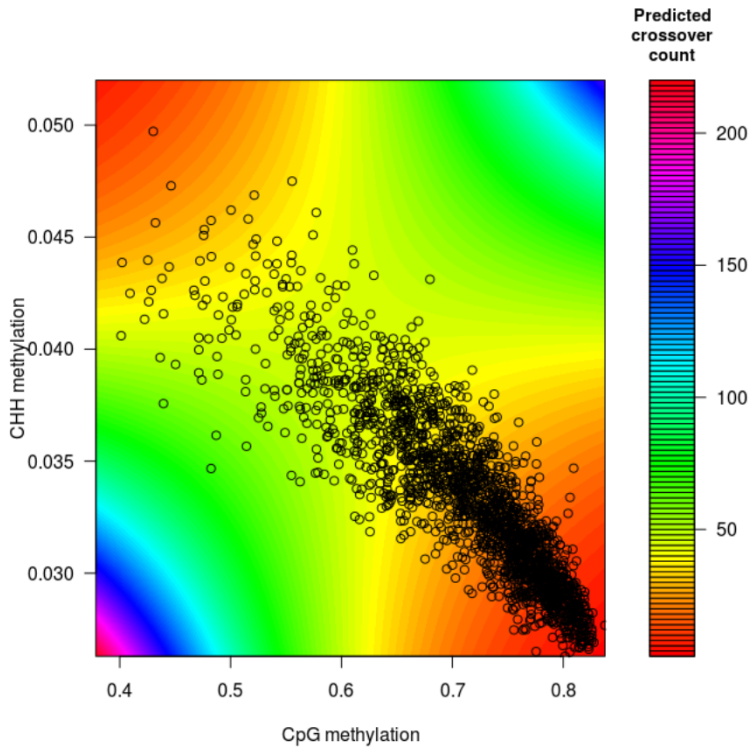
**Figure S7.** QQ-plot of  $\chi^2$  statistic testing number of crossovers across the US-NAM families proportional to the number of RILs in each family. The red line gives the threshold defined by the Bonferroni-corrected type I error rate,  $\alpha=0.05$



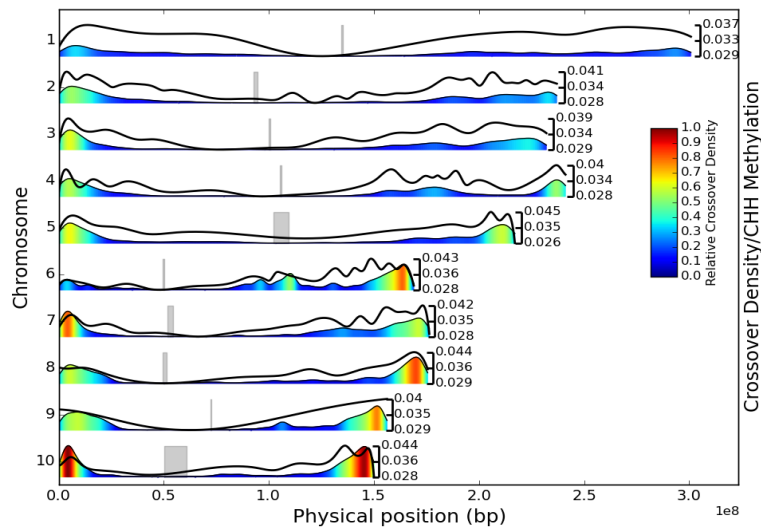
**Figure S8.** Pairwise scatter plots of the variables using in the 1-Mb model of crossover density in US-NAM



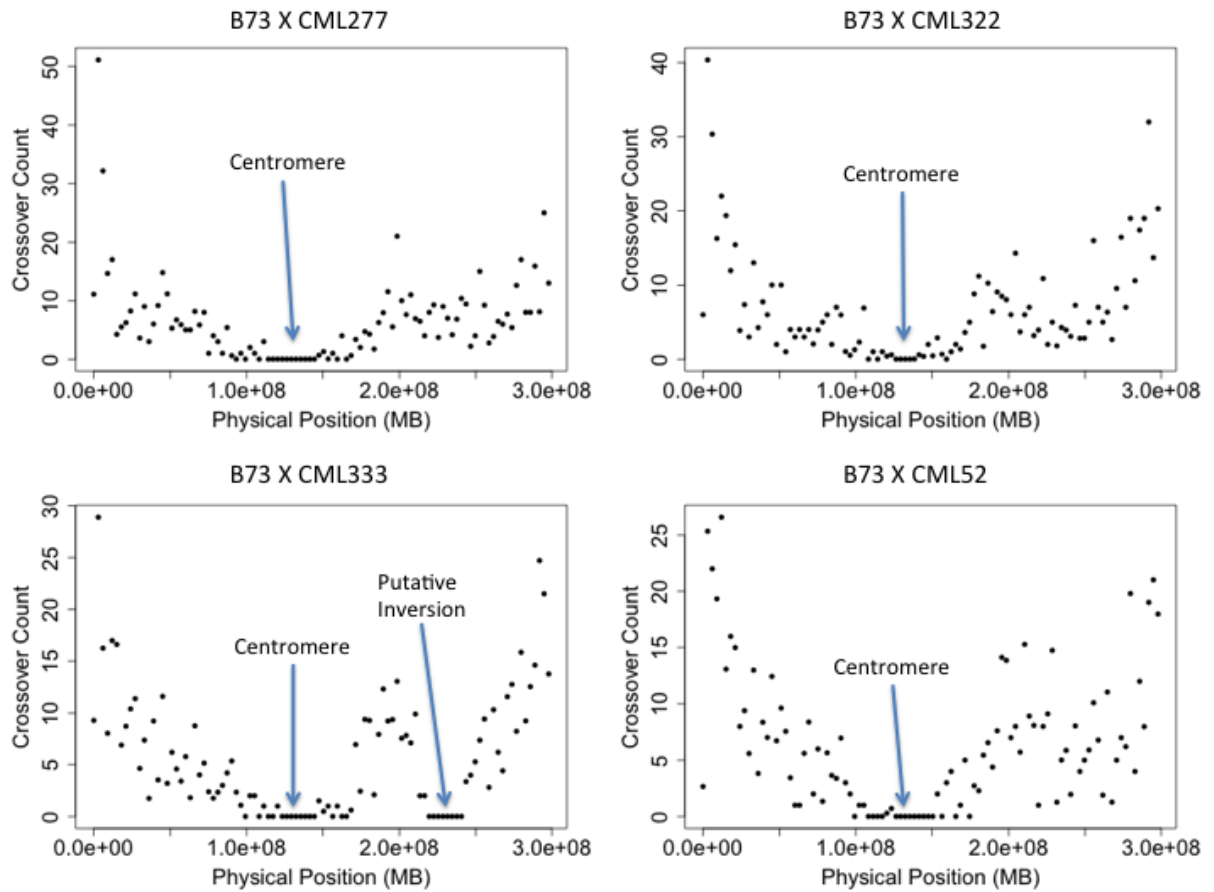
**Figure S9.** Genome-wide crossover density in US-NAM and its association with CHG methylation. Kernel density estimates of crossover density are shown by both height and color, drawn relative to the maximum density across all chromosomes, and black lines give the relative frequency of methylated CHGs, with scales given on the right side. Grey boxes indicate the locations of the centromeres.



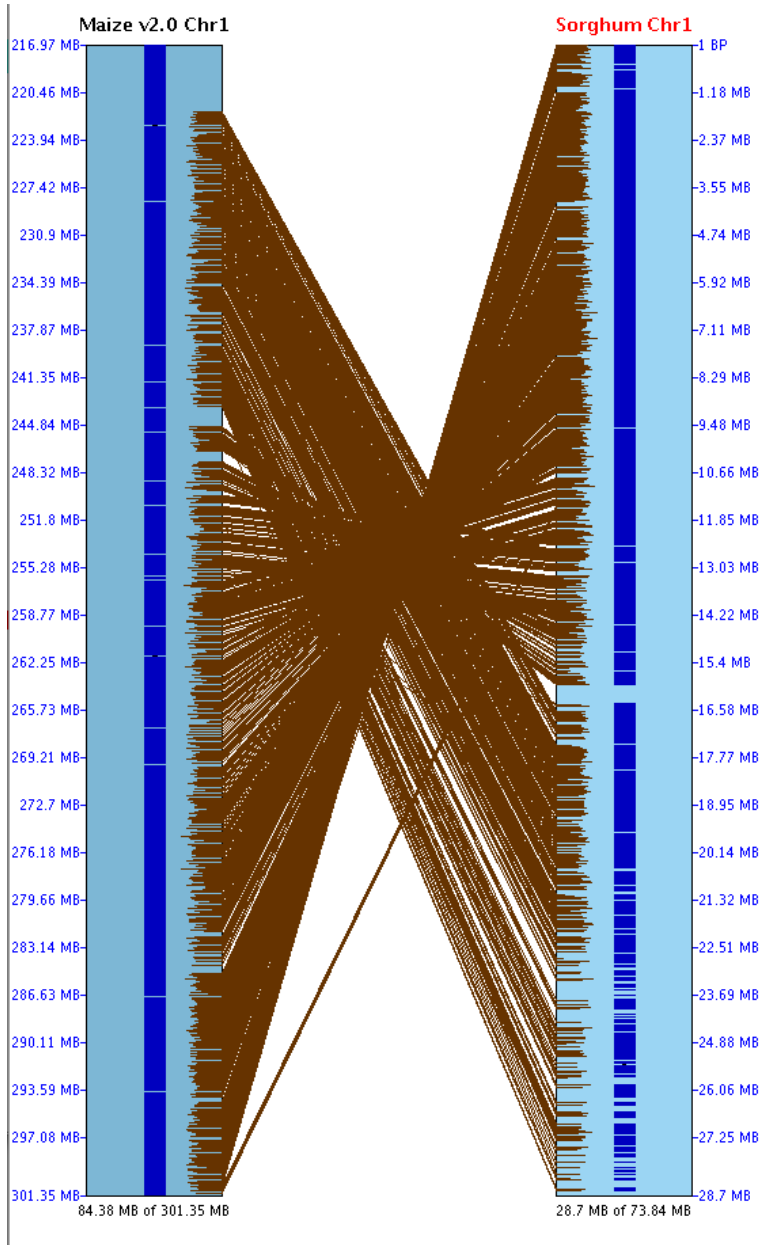
**Figure S10.** The predicted number of crossovers in 1Mb windows with varying CpG and CHH methylation according to the linear model. Other predictors are set at their mean values. Points show the observed values.



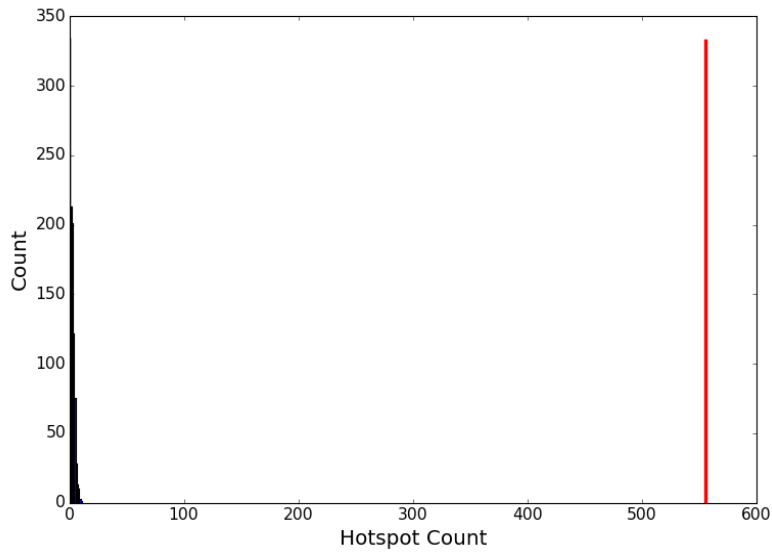
**Figure S11.** Genome-wide crossover density in US-NAM and its association with CHH methylation



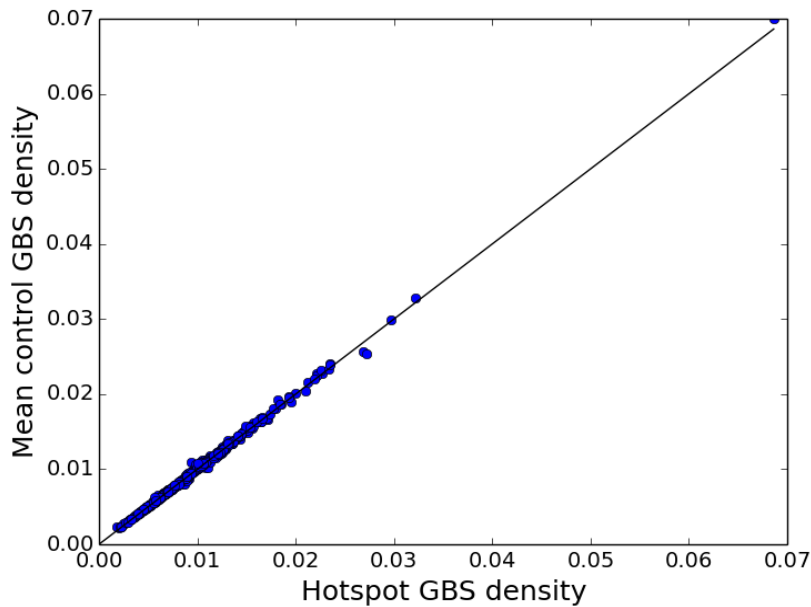
**Figure S12.** Crossover density in 4 US-NAM families on chromosome 1. In B73 x CML333 (Z007) there is a large non-centromeric region with no crossovers from 217.9 to 245.5 MB.



**Figure S13.** The inversion of the B73 chromosome 1 segment 217.9-245.48 Mb relative to sorghum chromosome 1. In B73 x CML333 this region contains no crossovers.

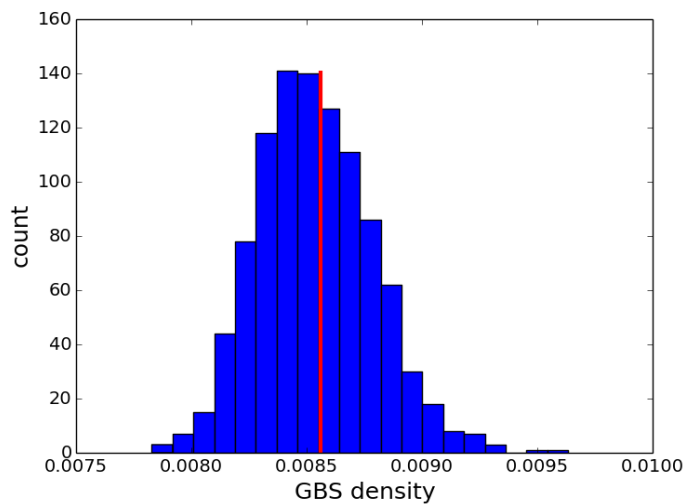


**Figure S14.** The total number of hotspots found using the smoothing spline-based method (red line) in the actual data compared to the number found in 1,000 simulations from a null distribution with the same Mb-scale pattern of recombination but without hotspots. Note that out of the total number of hotspots found, only 410 were used for further comparison in order to assure that each had at least 500 comparison controls.

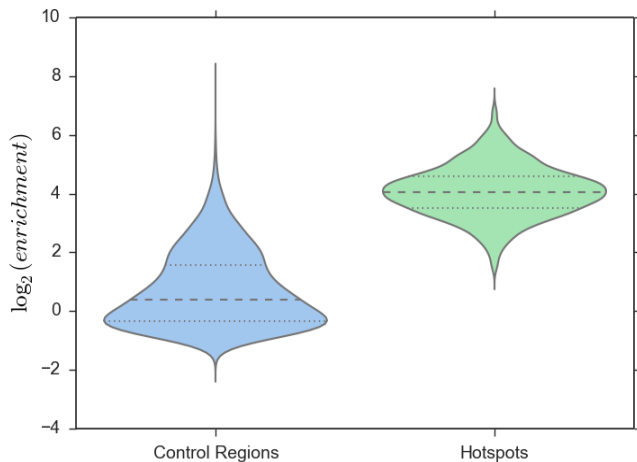


**Figure S15.** Relationship between the mean GBS density in controls and the GBS density of the corresponding hotspots.

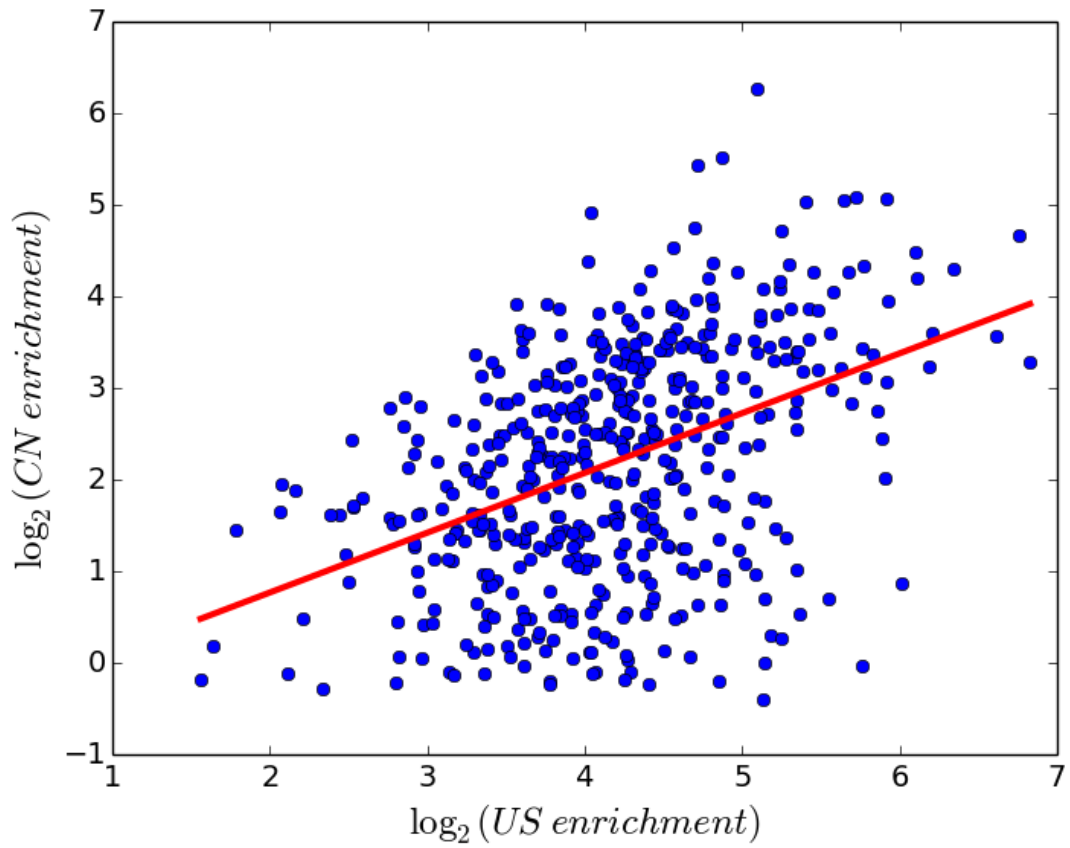




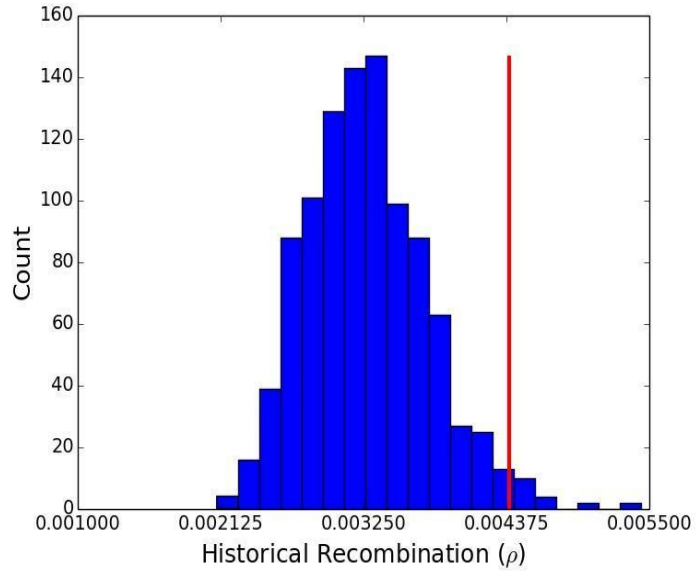
**Figure S16.** Mean GBS density of the hotspots (red line) compared to the range of mean GBS densities (blue histogram) of permuted controls during a permutation test.



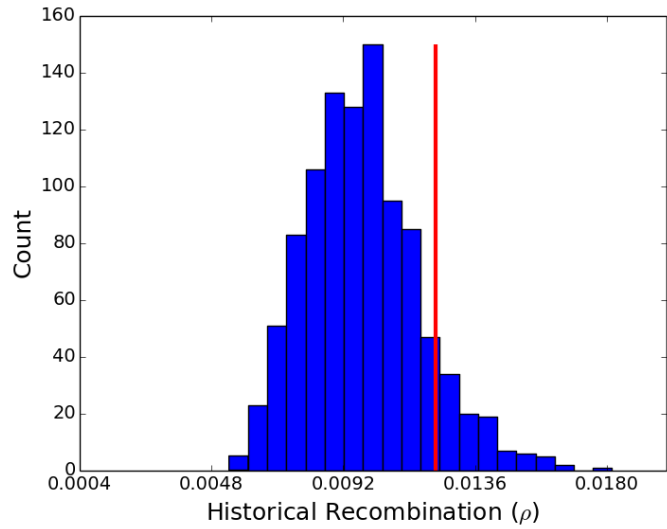
**Figure S17.** Estimated crossover enrichment ranges in controls and hotspots.



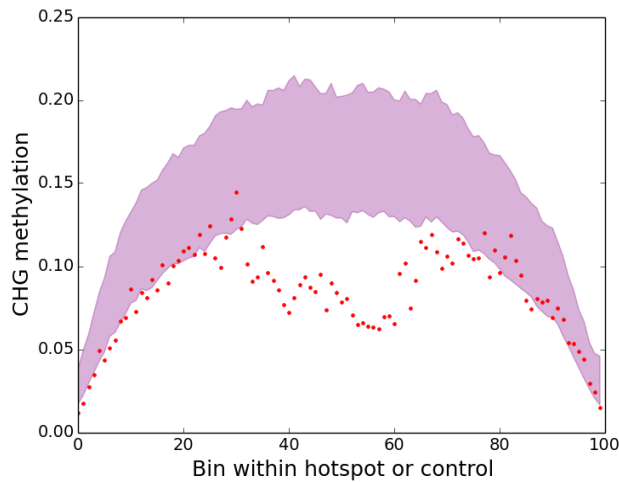
**Figure S18.** The relationship of estimated crossover enrichment in US-NAM hotspots to the crossover enrichment estimated in same CN-NAM regions



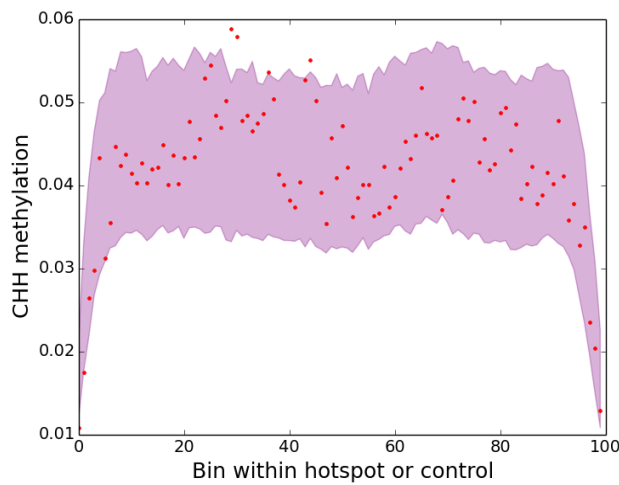
**Figure S19.** Mean historical recombination rate within maize landraces in hotspots compared to controls.



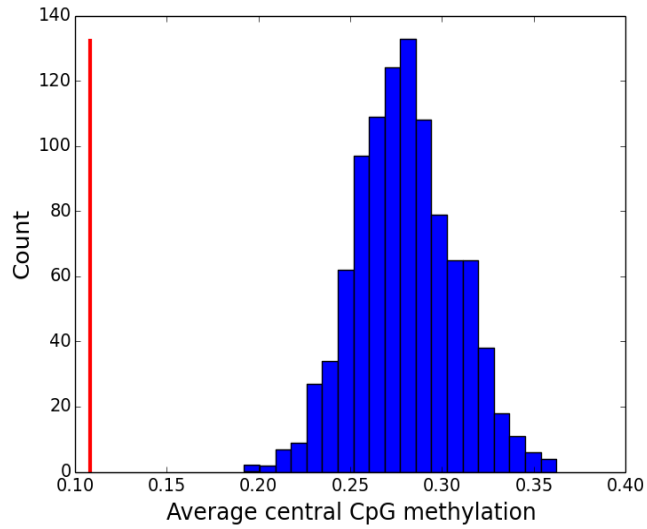
**Figure S20.** Mean historical recombination rate within teosintes in hotspots compared to controls.



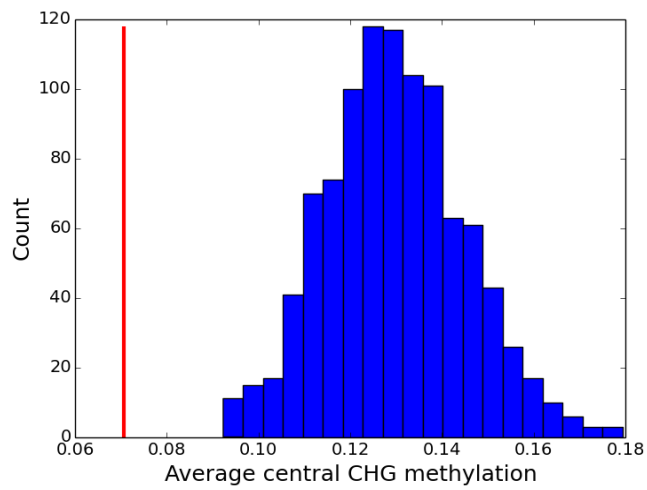
**Figure S21.** 95% CI for mean CHG methylation in 100 bins across controls (red shaded region) compared to 100 bins in hotspots (red dots). Each hotspot or control was divided into 100 even-sized bins, and the amount of methylation was averaged over all regions for each bin.



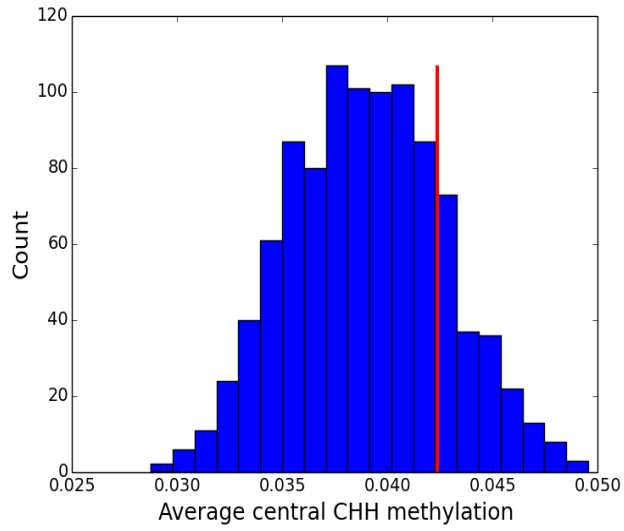
**Figure S22.** 95% CI for mean CHH methylation in 100 bins across controls (red shaded region) compared to 100 bins in hotspots (red dots).



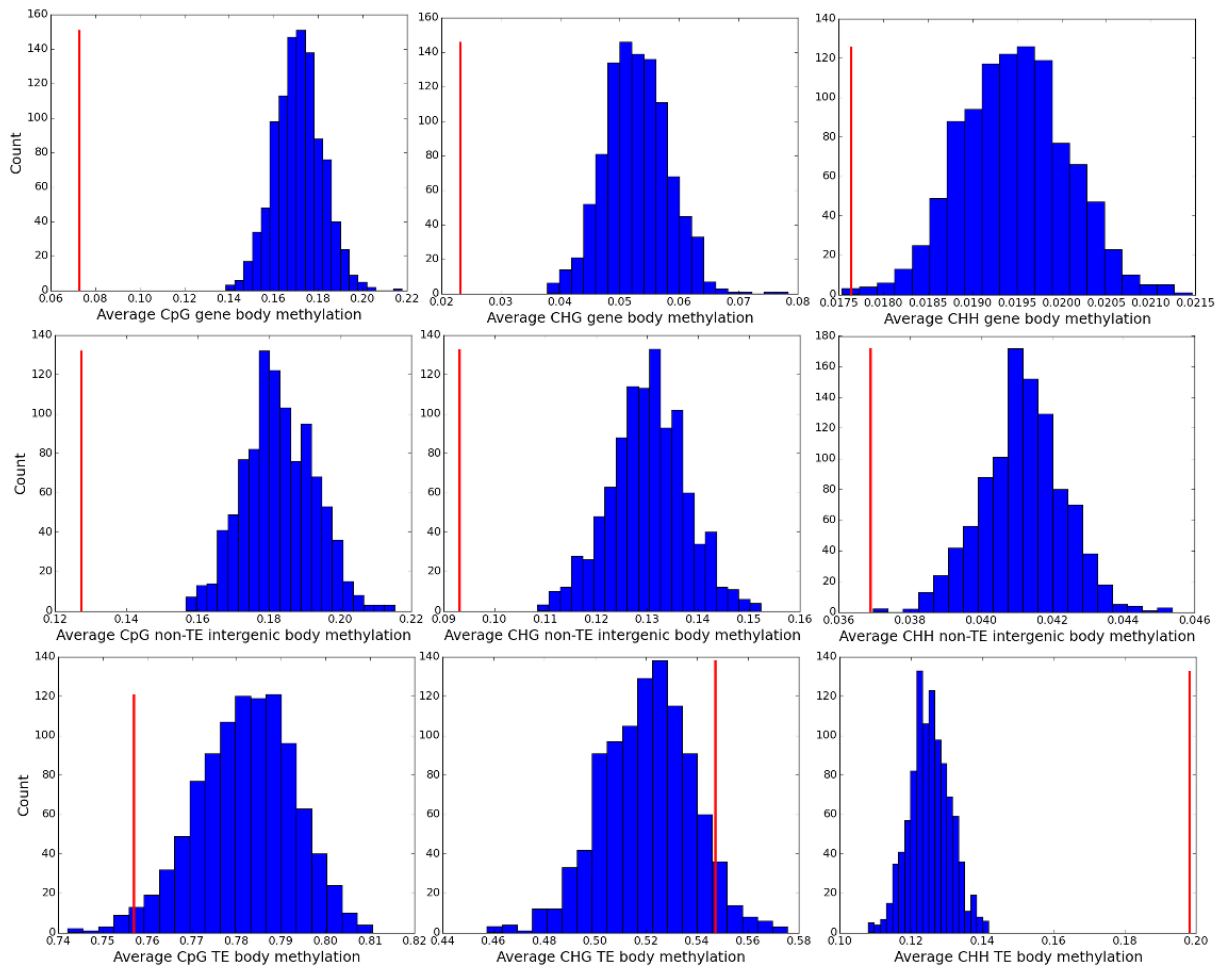
**Figure S23.** Mean central CpG methylation of hotspots compared to controls that have mean GBS depth greater than or equal to their comparison hotspot.



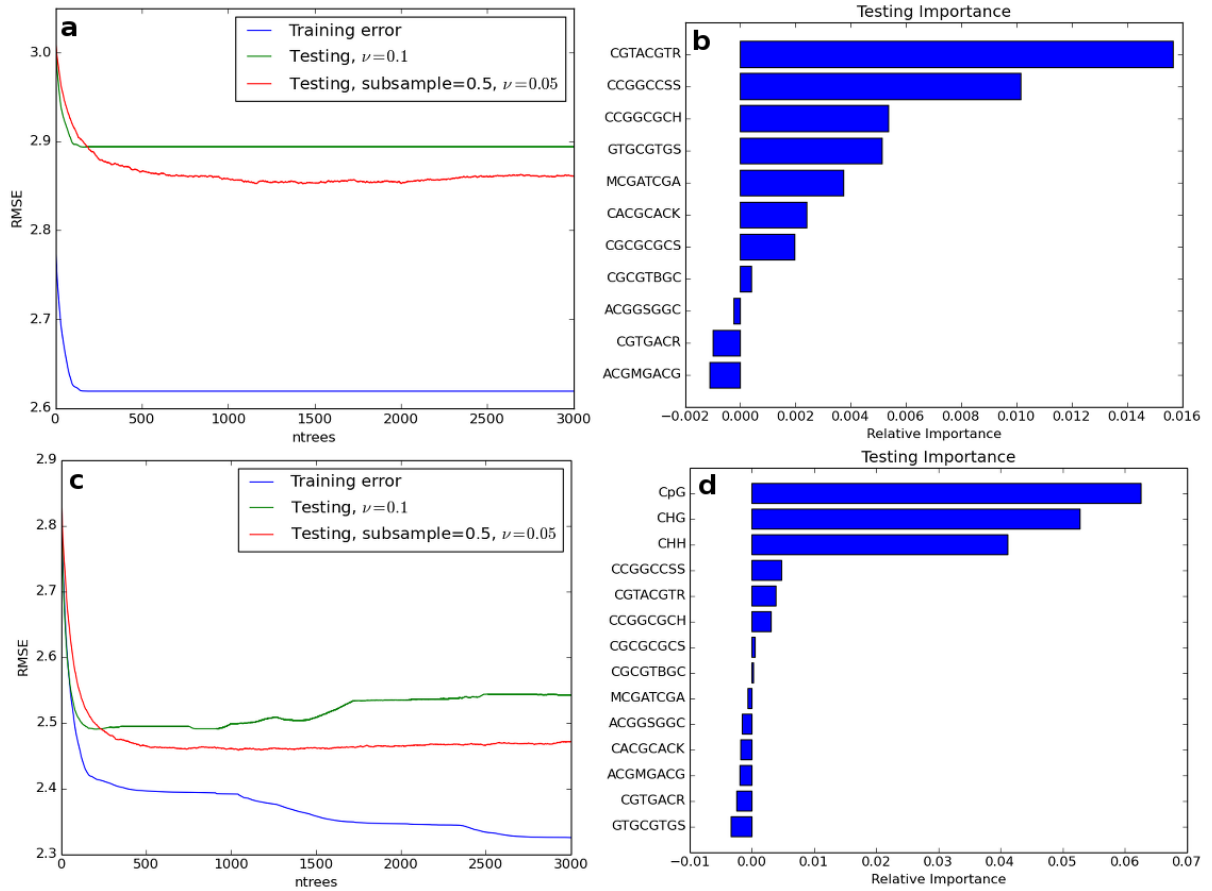
**Figure S24.** Mean central CHG methylation of hotspots compared to controls that have mean GBS depth greater than or equal to their comparison hotspot.



**Figure S25.** Mean central CHH methylation of hotspots compared to controls that have mean GBS depth greater than or equal to their comparison hotspot.

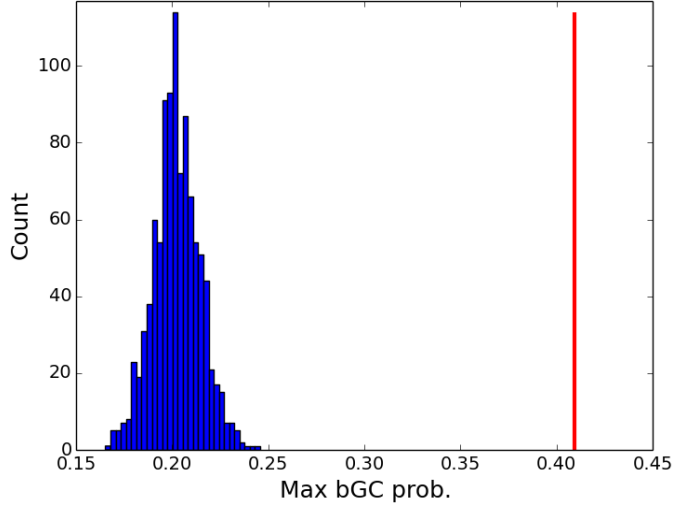


**Figure S26.** Mean methylation of hotspots compared to controls, separated by genomic context (gene body, transposable element body, non-TE intergenic region).

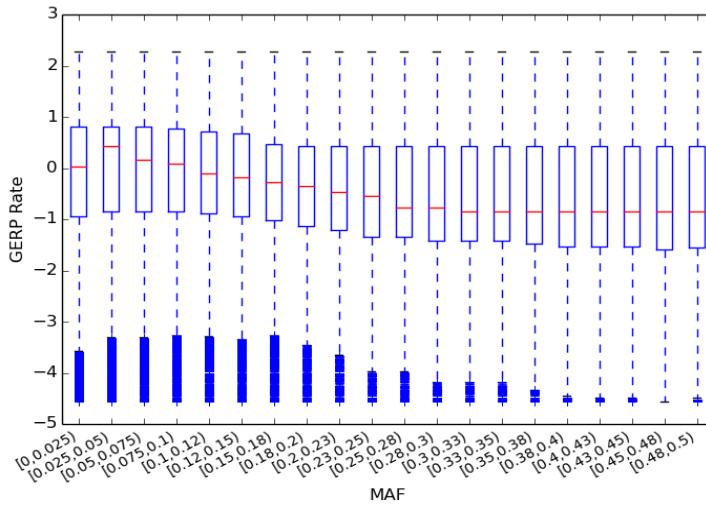


**Figure S27.** Results from using machine learning to predict the number of crossovers in 30kb intervals using only motifs significantly enriched in hotspots (a and b) and using the motifs in addition to the amount of CpG, CHG, and CHH methylation (c and d).

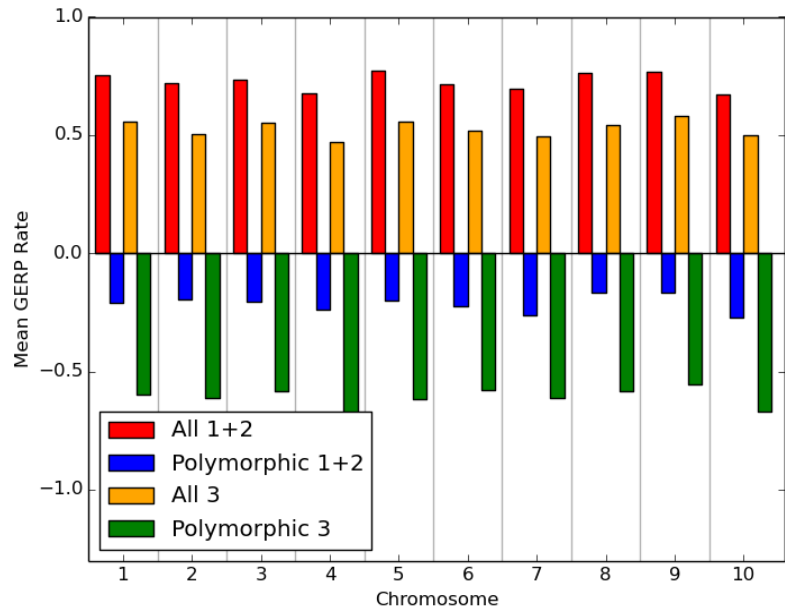




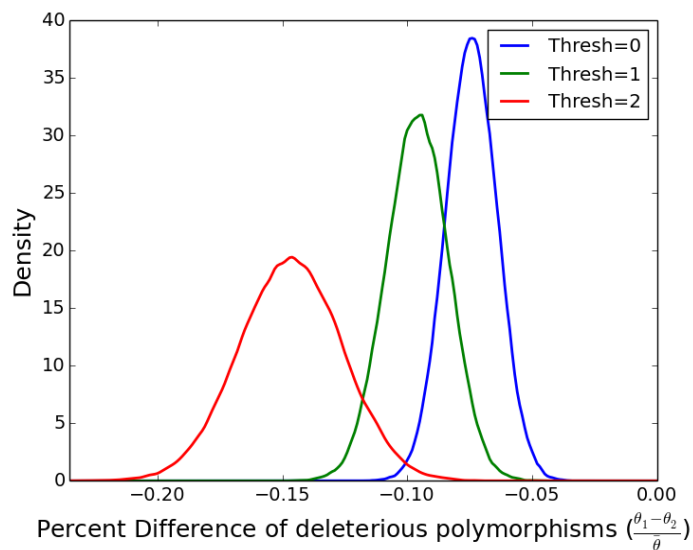
**Figure S28.** Mean maximum posterior bGC probability in hotspots compared to controls.



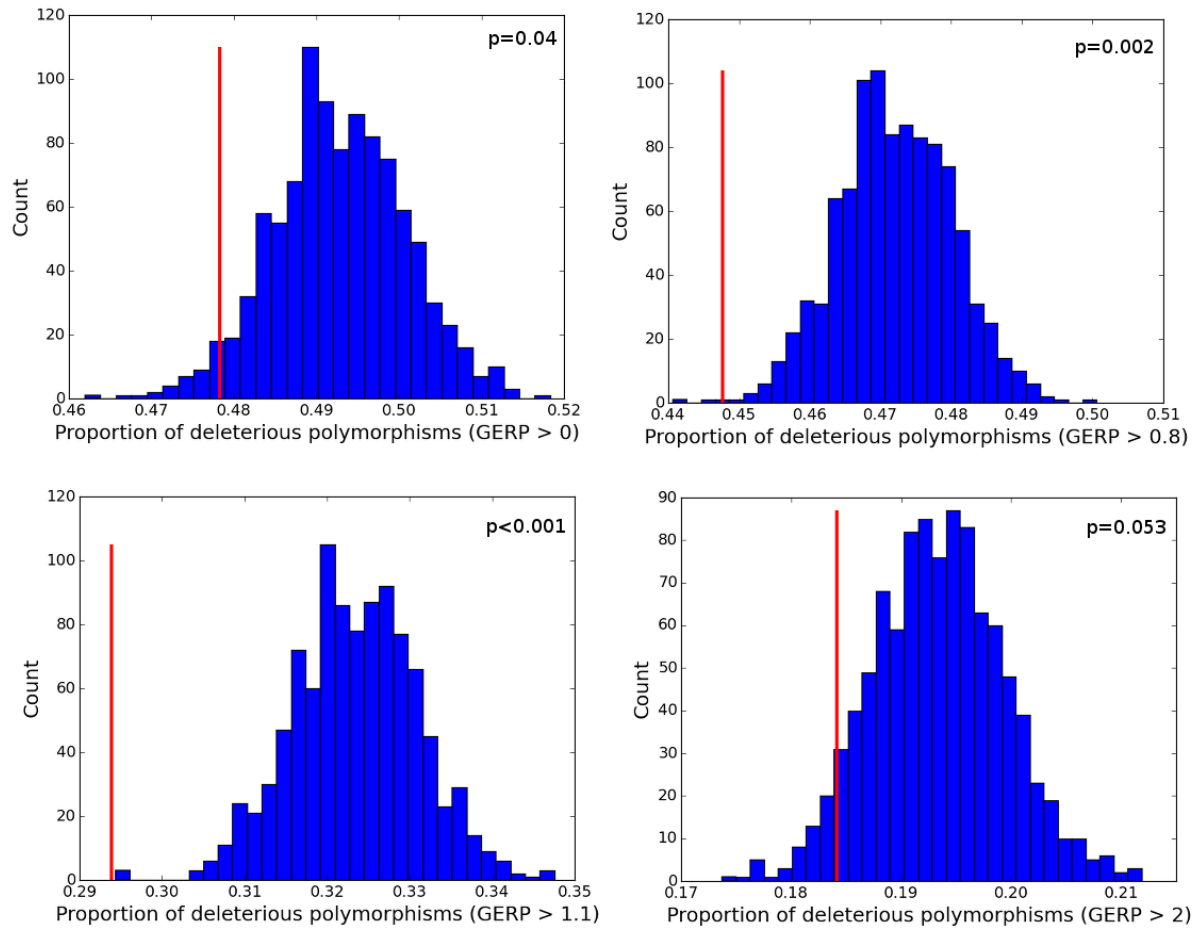
**Figure S29.** GERP rates across the genome, grouped by the minor allele frequency. In this study we removed SNPs with MAFs below 0.025 in order to guard against the effects of genotyping error.



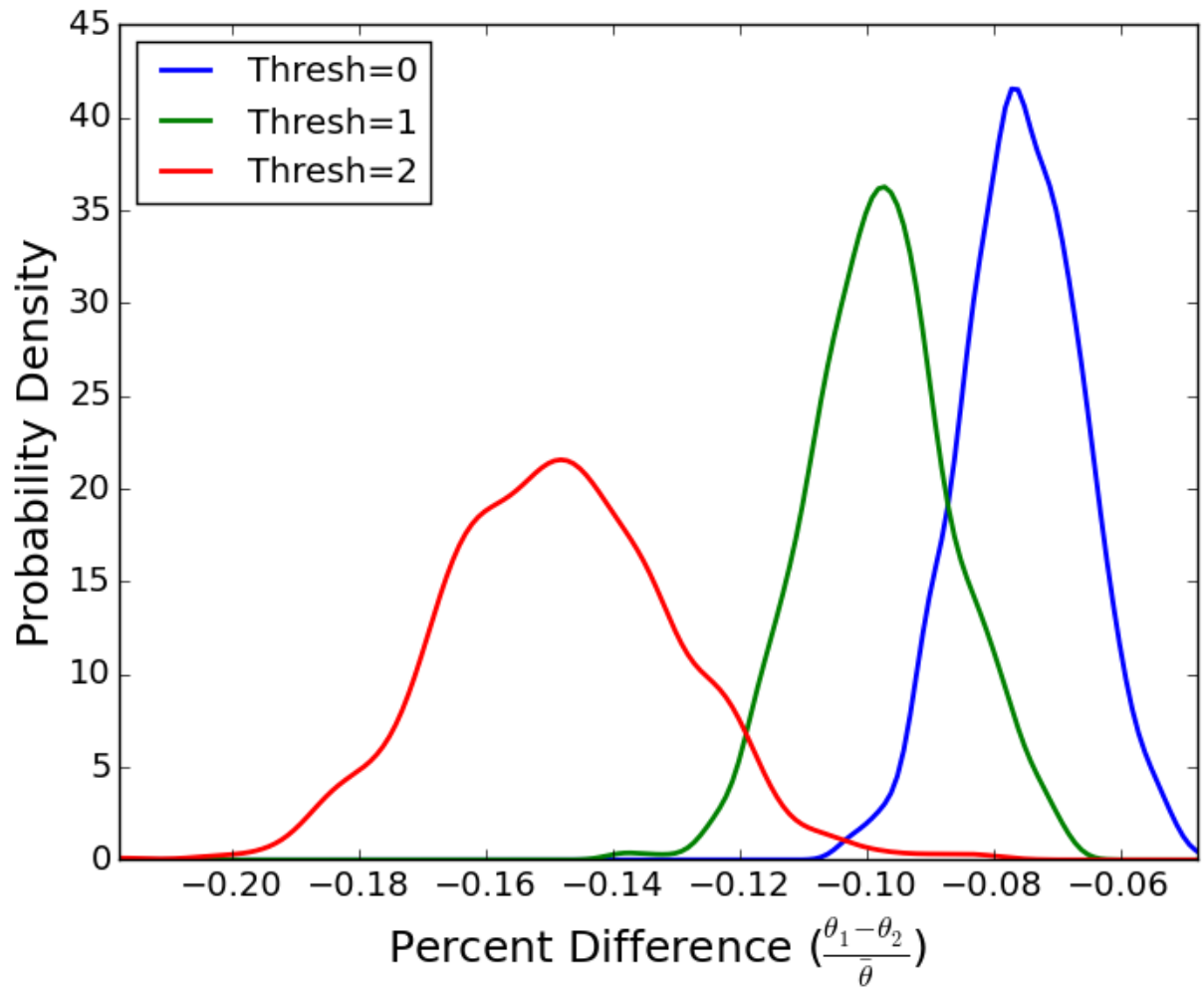
**Figure S30.** The mean GERP rates within annotated codons across the genome. Bars give the mean rates at all first and second codon positions or all third codon positions (red and orange bars) or the mean rates at those positions with known segregating polymorphisms (blue and green bars).



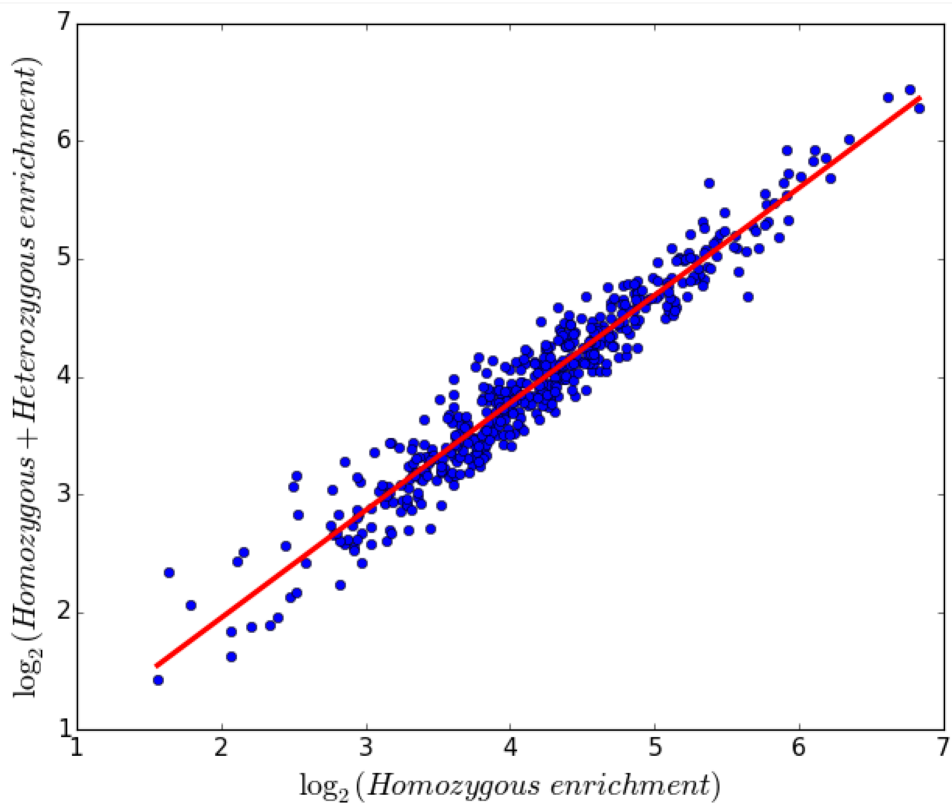
**Figure S31.** Percent difference between the proportion of polymorphisms at sites with GERP above a threshold in hotspots ( $\theta_1$ ) vs. the rest of the genome ( $\theta_2$ ).



**Figure S32.** The proportion of deleterious polymorphisms, as measured at varying GERP thresholds, of hotspots compared to controls after limiting control regions to those with GC content at least as high as the comparison hotspot. One-tailed permutation p-values are shown.



**Figure S33.** Kernel density estimate of the percent difference between the proportion of polymorphisms at sites with GERP above a threshold in hotspots ( $\theta_1$ ) vs. the rest of the genome ( $\theta_2$ ) based on permutations of random sites from the genome in which the proportion of each reference base was constrained to be the same as the proportions within the hotspots.



**Figure S34.** Relationship of estimated crossover enrichment in recombination hotspots between datasets where heterozygous crossovers are included and excluded.