

1 **Cassava HapMap: Masking deleterious mutations in a clonal crop species**

2

3 **Punna Ramu<sup>1\*</sup>, Williams Esuma<sup>2</sup>, Robert Kawuki<sup>2</sup>, Ismail Y Rabbi<sup>3</sup>, Chiedozi Egesi<sup>3,4,5</sup>,**  
4 **Jessen V Bredeson<sup>6</sup>, Rebecca S Bart<sup>7</sup>, Janu Verma<sup>1</sup>, Edward S Buckler<sup>1,8</sup>, Fei Lu<sup>1\*</sup>**

5

6 <sup>1</sup>Institute of Genomic Diversity, Cornell University, Ithaca, NY, USA.

7 <sup>2</sup>National Crops Resources Research Institute (NaCRRI), Kampala, Uganda.

8 <sup>3</sup>International Institute of Tropical Agriculture (IITA), Ibadan, Nigeria.

9 <sup>4</sup>National Root Crops Research Institute (NRCRI), Umudike, Nigeria.

10 <sup>5</sup>International Programs, College of Agriculture and Life Sciences, Cornell University,  
11 Ithaca, NY, USA.

12 <sup>6</sup>Department of Molecular and Cell Biology, University of California, Berkeley, CA, USA.

13 <sup>7</sup>Donald Danforth Plant Science Center, St. Louis, MO, USA.

14 <sup>8</sup>US Department of Agriculture – Agriculture Research Service (USDA-ARS).

15 Correspondence should be addressed to P.R. ([rp444@cornell.edu](mailto:rp444@cornell.edu)) or F.L.

16 ([fl262@cornell.edu](mailto:fl262@cornell.edu))

17 **Cassava (*Manihot esculenta* Crantz) is an important staple food crop in Africa and**  
18 **South America, however, ubiquitous deleterious mutations may severely reduce its**  
19 **fitness. To evaluate these deleterious mutations in the cassava genome, we**  
20 **constructed a cassava haplotype map using deep sequencing from 241 diverse**  
21 **accessions and identified over 28 million segregating variants. We found that, 1) while**  
22 **domestication modified starch and ketone metabolism pathways for human**  
23 **consumption, the concomitant bottleneck and clonal propagation resulted in a large**  
24 **proportion of fixed deleterious amino acid changes, raised the number of deleterious**  
25 **mutations by 26%, and shifted the mutational burden towards common variants; 2)**  
26 **deleterious mutations are ineffectively purged due to limited recombination in**  
27 **cassava genome; 3) recent breeding efforts maintained the yield by masking the most**  
28 **damaging recessive mutations in the heterozygous state, but unable to purge the**  
29 **mutation burden, which should be a key target for future cassava breeding.**

30

31 Cassava is the third most consumed carbohydrate source for millions of people in  
32 tropics, after rice and maize<sup>1</sup>. Even though cassava was domesticated in Latin America<sup>2</sup>,  
33 it has spread widely and become a major staple crop in Africa. Cassava stores starch in  
34 underground storage roots, which remain fresh until harvest. Cassava is a highly  
35 heterozygous species. Although its wild progenitor, *M. esculenta* ssp. *falbellifolia*,  
36 reproduces by seed<sup>3</sup>, it is particularly worth noted that cultivated cassava is almost  
37 exclusively clonally propagated via stem cutting, in which a single individual contributes  
38 its entire genome to its offspring<sup>4</sup>. The limited number of recombination events in such  
39 vegetatively propagated crops results in a potential accumulation of deleterious  
40 mutations across the genome<sup>5</sup>. Thus, mutation burden in cassava is expected to be  
41 more severe than in sexually propagated species. Deleterious mutations are considered  
42 to be at the heart of inbreeding depression<sup>6</sup>. Inbreeding depression is extremely severe,  
43 even in elite cassava accessions, where a single generation of inbreeding results in >60%

44 reduction in fresh root yield<sup>7,8</sup>. In this study, we aimed to identify deleterious mutations  
45 in cassava populations, which in turn can help accelerate cassava breeding by allowing  
46 breeders to purge deleterious mutations more efficiently.

47 We conducted a comprehensive characterization of genetic variation by whole genome  
48 sequencing (WGS) of 241 cassava accessions, including 203 elite breeding accessions (*M.*  
49 *esculenta* Crantz), 16 close relatives (*M. esculenta* ssp. *flabellifolia*, *M. esculenta* ssp.  
50 *peruviana*) of modern cultivars<sup>2,9</sup>, 11 hybrid/tree cassava accessions, and 11 more  
51 divergent wild relatives (*M. glaziovii* and others) (**Supplementary Fig. 1** and  
52 **Supplementary Table 1**). Samples included 54 accessions from an initial haplotype map I  
53 (HapMapI) study<sup>10</sup>. Wild *M. glaziovii* has been used extensively in cassava breeding  
54 programs to transfer disease resistance alleles to cultivated cassava (e.g., Amani  
55 Breeding program)<sup>8</sup>. On average, more than 30x coverage sequences were generated  
56 for each accession. The 518.5 Mb cassava genome (v6.1) has roughly 51% repetitive  
57 elements with several common recent retrotransposons<sup>10</sup>. To exclude misalignment and  
58 ensure high quality of variant calling, repeat sequences were pre-filtered using repeat  
59 bait (**Supplementary Fig. 2**) and the remaining sequences were aligned against the  
60 cassava reference genome v6.1<sup>10,11</sup>. Variants from low copy regions of the genome were  
61 identified to develop the cassava haplotype map II (HapMapII) with 27.8 million variants  
62 (25.9 million SNPs and 1.9 million indels) and with a low error rate of 0.01%, which is the  
63 proportion of segregating sites in the reference accession (**Supplementary Fig. 3**). The  
64 correlation between read depth and proportion of heterozygotes of SNPs is extremely  
65 low ( $r^2 = 6E-05$ , **Supplementary Fig. 4**). Cultivated cassava exhibited 9.94 million variants  
66 (**Supplementary Table 2**), of which nearly 50% were found to be rare (<5% minor allele  
67 frequency (MAF)) (**Supplementary Table 2 and Supplementary Fig. 5**). Haplotypes were  
68 phased and missing genotypes were imputed with high accuracy using BEAGLE v4.1<sup>12</sup>  
69 (accuracy  $r^2 = 0.966$ ) (**Supplementary Fig. 6**). Linkage disequilibrium was as low as in  
70 maize<sup>13</sup> and decayed to an average  $r^2 = 0.1$  in 3,000 bp (**Supplementary Fig. 7**).

71 Cultivated cassava presented lower nucleotide diversity ( $\pi = 0.0036$ ) compared with its  
72 progenitors (*M. esc.* ssp. *flabellifolia*,  $\pi = 0.0051$ ). In addition, a close relationship  
73 between the two species was observed from phylogenetic analysis (**Supplementary Fig.**  
74 **8**). Both lines of evidence support the hypothesis that cultivated cassava was  
75 domesticated from *M. esc.* ssp. *flabellifolia*<sup>2,9,10</sup>. To evaluate population differentiation  
76 of cassava, a principal component (PC) analysis was performed and showed substantial  
77 differentiation among all cassava species and hybrids (**Fig. 1a**), where cultivated cassava  
78 showed moderate genetic differentiation from its progenitors ( $F_{st}$ : 0.16), and high  
79 genetic differentiation from tree cassava ( $F_{st}$ : 0.32) and wild relatives ( $F_{st}$ : 0.44)  
80 (**Supplementary Table 2 and Supplementary Figs. 9 and 10**). However, PC analysis  
81 showed very little differentiation among cultivated cassava (**Fig. 1b**), where geographic  
82 subpopulations of cultivated cassava presented surprisingly low value of  $F_{st}$  among  
83 themselves (0.01-0.05) despite the fact that these subpopulations were sampled from  
84 different continents (**Supplementary Table 2**). This suggests that despite clonal

85 propagation, there has been enough crossing to keep cultivated cassava in one breeding  
86 pool.

87 Sequence conservation is a powerful tool to discover functional variation<sup>14,15</sup>. We  
88 identified deleterious mutations by utilizing genomic evolution and amino acid  
89 conservation modeling. The cassava genome was aligned to seven species in the  
90 Malpighiales clade to identify evolutionarily constrained regions of cassava genome.  
91 Based on genomic evolutionary rate profiling (GERP)<sup>16</sup> score, nearly 104-Mb of the  
92 genome (20%) of cassava was constrained (GERP score > 0) (**Supplementary Fig. 11**).  
93 The evolutionarily constrained genome of cassava (104 Mb) is comparable to maize (111  
94 Mb)<sup>17</sup> in size, but less than humans (214 Mb)<sup>16</sup> and more than *Drosophila* (88 Mb)<sup>18</sup>.  
95 GERP profiling also identified remarkably asymmetric distribution of constrained  
96 sequence at the chromosome scale (**Supplementary Fig. 12**). In addition to the  
97 constraint estimation at the DNA level, consequences of mutation on amino acids in  
98 proteins were assessed using Sorting Intolerant From Tolerant (SIFT) program<sup>19</sup>. Nearly  
99 3.0% of coding SNPs in cultivated cassava were non-synonymous mutations  
100 (**Supplementary Table 2 and Supplementary Fig. 13**), of which 19.3% (57,952) were  
101 putatively deleterious (SIFT < 0.05). As the strength of functional prediction methods  
102 varies<sup>14</sup>, we combined SIFT (< 0.05) and GERP (> 2) to obtain a more conservative set of  
103 22,495 deleterious mutations (**Supplementary Fig. 14**).

104 To estimate the individual mutation burden, we used rubber (*Hevea brasiliensis*), which  
105 diverged from the cassava lineage 27 million years ago<sup>10</sup>, as an out-group to identify  
106 derived deleterious alleles in cassava. First, we focused on the fixed deleterious  
107 mutations. The derived allele frequency (DAF) spectrum shows that cassava (5%, **Fig. 2**)  
108 appears to have more fixed deleterious mutations than maize (3.2%, DAF > 0.8)<sup>20</sup> when  
109 compared at the same threshold (SIFT < 0.05). Across cultivated cassava there were 150  
110 fixed deleterious mutations. These deleterious mutations cannot be purged through  
111 standard breeding which relies on recombination of segregating alleles, but these fixed  
112 deleterious mutations are the potential targets for genome editing<sup>21</sup>. Together with the  
113 other 22,345 segregating deleterious mutations, the mutation burden in cassava was  
114 substantial. Given the several millennia of breeding in the species, why are these  
115 deleterious mutations still in cultivated cassava and how were breeders managing  
116 them? We evaluated the effects of recombination, selection, and drift, as the main  
117 processes controlling the distribution of deleterious mutations in the genome.

118 Recombination is an essential process to purge deleterious mutations from genome<sup>22</sup>. In  
119 vegetatively propagated species like cassava, recombination is expected be less efficient  
120 in purging deleterious mutations. This hypothesis was supported by a weak correlation  
121 between recombination rate and distribution of deleterious mutations ( $r = -0.065$ ,  $P =$   
122  $0.13$ , **Fig. 3a**). Deleterious mutation were nearly uniformly spread across the cassava  
123 genome (**Fig. 3b and Supplementary Fig. 15**), rather than being concentrated in low  
124 recombination regions as in human<sup>23</sup>, fruit fly<sup>24</sup>, and maize<sup>17</sup>. Thus, recombination,

125 which is presumably rare in a clonally propagated crop, does not effectively purge  
126 mutation burden in cassava.

127 Domestication is important in evolution and improvement of crop species. The major  
128 domestication trait of cassava is the large carbohydrate rich storage root. Cultivated  
129 cassava has 5-6 times higher starch content than its progenitor<sup>3</sup>. Another domestication  
130 trait is the reduced cyanide content in roots<sup>3</sup>. Every tissue of cassava contains  
131 cyanogenic glucosides<sup>25</sup>. Ketones, cyanohydrin, and hydrogen cyanide are the key toxic  
132 compounds formed upon degradation of cyanogenic glucosides<sup>25,26</sup>. These toxic  
133 compounds have to be eliminated before consumption. To identify the genomic regions  
134 under selection during the domestication, a likelihood method (the cross-population  
135 composite likelihood ratio, XP-CLR)<sup>27</sup> was used to scan the genome in Latin American  
136 accessions and the progenitor *M. esculenta* ssp. *flabellifolia*. We identified 203 selective  
137 sweeps containing 427 genes in Latin American accessions (**Supplementary Fig. 16a**).  
138 Genes in these sweep regions were enriched for starch and sucrose synthesis (3.8-fold  
139 enrichment; FDR =  $7.2 \times 10^{-03}$ ) and cellular ketone metabolism (3.4-fold enrichment; FDR  
140 =  $5.3 \times 10^{-03}$ ) (**Supplementary Fig. 16b**). The results suggest that selection during  
141 domestication increased production of carbohydrates and reduced cyanogenic glucoside  
142 in cassava. Likewise, selection signatures of recent bottleneck event in African cassava  
143 accessions were also evaluated. A total of 244 selective sweeps were identified  
144 containing 416 genes. These genes were enriched for serine family amino acid  
145 metabolism (4.2-fold enrichment, FDR =  $2.1 \times 10^{-06}$ ) and cellular response to stress (1.3-  
146 fold enrichment, FDR =  $4.9 \times 10^{-06}$ , **Supplementary Fig. 17**). Since L-Serine is involved in  
147 the plant response to biotic and abiotic stresses<sup>28,29</sup>, together with the functional  
148 enrichment of cellular response to stress, it may reflect that disease resistance  
149 accessions were selected in recent breeding program in Africa<sup>8</sup>.

150 How was the genetic burden shaped in the selective sweeps? We found that Latin  
151 American accessions showed 25% less ( $P = 0.009$ , **Fig. 4a**) deleterious mutations than  
152 progenitors in sweep regions. Similarly, African accessions exhibited a 35% drop ( $P = 2.1$   
153  $\times 10^{-07}$ , **Fig. 4b**) in sweeps compared to Latin American accessions. In addition to the  
154 comparison between populations, significant reductions of deleterious mutations were  
155 observed within population by comparing sweep regions and the rest of the genome.  
156 For example, selective sweeps presented 44% depletion ( $P = 9.7 \times 10^{-12}$ , **Fig. 4c**) of  
157 deleterious mutations in Latin American accessions and 41% reduction ( $P = 8.7 \times 10^{-130}$ ,  
158 **Fig. 4d**) in African accessions. This implies that haplotypes containing fewer deleterious  
159 alleles were favored during selection.

160 However, drift after domestication played a more important role in affecting mutation  
161 burden in cassava. Although Latin American accessions and African accessions had a  
162 similar number of deleterious mutations ( $P = 0.42$ , **Fig. 5a**), they presented a prominent  
163 increase of total burden by 26% ( $P = 9.1 \times 10^{-09}$ , **Fig. 5a**) when compared with  
164 progenitors, and shifted the mutation burden towards common deleterious variants  
165 (**Supplementary Fig. 18**). The increase of deleterious mutations during domestication

166 was also found in dog<sup>30</sup>. The results suggest that the severe bottleneck of domestication  
167 and shift from sexual reproduction to clonal propagation resulted in a rapid  
168 accumulation of deleterious mutations in cultivated cassava.

169 How have the breeders been able to maintain yield, given the substantial growth of  
170 mutation burden in cultivated cassava? This became apparent when the homozygous  
171 deleterious mutations and heterozygous deleterious mutations were compared.  
172 Relative to *M. esculenta* ssp. *flabellifolia*, the homozygous mutation burden  
173 substantially decreased by 23% ( $P = 7 \times 10^{-03}$ , **Fig. 5b**) in cultivated accessions regardless  
174 of the elevated frequency of deleterious alleles (**Supplementary Fig. 18**), while the  
175 heterozygous mutation burden remarkably increased by 96% ( $P = 8.1 \times 10^{-07}$ , **Fig. 5c**),  
176 despite the reduced genetic diversity in cultivated cassava ( $\pi = 0.0036$ ) relative to  
177 progenitors ( $\pi = 0.0051$ ). In addition to the comparison between cultivated cassava and  
178 progenitors, we also compared observed and expected mutation burden under the  
179 assumption of Hardy–Weinberg Equilibrium (HWE) within cultivated cassava (**Online**  
180 **Methods**). Although HWE was probably never reached in the breeding pool, the relative  
181 depletion of homozygous mutation burden and excess of heterozygous mutation  
182 burden would not be seen unless it was selected and maintained. The results from  
183 bootstrap resampling (10,000 times) showed that the observed homozygous mutation  
184 burden was less than the expected (Latin American cassava: 5.6% decrease,  $P = 0$ ;  
185 African cassava: 10.3% decrease,  $P = 0$ , **Fig. 5d**), and the observed heterozygous  
186 mutation burden was more than expected (Latin American cassava: 3.5% increase,  $P =$   
187  $1.5 \times 10^{-312}$ ; African cassava: 6.9% increase,  $P = 0$ , **Fig. 5e**), indicating a significant  
188 deviation from HWE expectation. These evidence suggest that breeders have been  
189 trying to manage the recessive deleterious mutations in the heterozygous state to mask  
190 their harmful effects.

191 Mutations with large homozygous effect are more likely to be recessive<sup>31</sup>. We found  
192 nearly 64.5% of deleterious mutations occurred only in the heterozygous state  
193 (**Supplementary Fig. 19**). Although the low allele frequency confines effective tests for  
194 the excess heterozygotes of these deleterious mutations, they are more likely to be  
195 strong deleterious mutations, resulting in the significant yield loss in the first generation  
196 of selfed cassava plants<sup>7,8</sup>. These mutations were in genes ( $n = 7,774$ ) mainly enriched  
197 for macromolecule catabolism and biosynthesis (**Supplementary Fig. 20a**). In contrast,  
198 the deleterious mutations existing predominantly in the homozygous state (proportion  
199 of homozygotes > 70%, **Supplementary Fig. 19**), were present in genes ( $n = 245$ )  
200 enriched for amine and ketone metabolism, as well as chemical and stimulus responses  
201 (**Supplementary Fig. 20b**). This difference suggests that the deleterious mutations  
202 primarily exhibited in the heterozygous state may have relatively large fitness  
203 consequences.

204 Cassava is a major staple crop feeding hundreds of millions of people. Using deep  
205 sequencing of a comprehensive and representative collection of 241 cassava accessions,  
206 we developed the HapMapII, a highly valuable resource for cassava genetic studies and

207 breeding. In this vegetatively propagated species, deleterious mutations have been  
208 accumulating rapidly due to the lack of recombination. The bottleneck event during  
209 domestication exacerbated the existing mutation burden in cassava. Breeding efforts  
210 successfully maintained the yield by selecting high fitness haplotypes at a few hundred  
211 loci and handling most damaging mutations in the heterozygous state. However,  
212 breeders were unable to purge the mutation burden due to limited recombination,  
213 instead they shielded deleterious mutations by increasing the heterozygosity while  
214 screening thousands of potential hybrids (**Supplementary Fig. 21**). In the short term,  
215 this practice for managing mutation burden may produce gains in yield. In the long run,  
216 however, a mutational meltdown may be triggered by new mutations, decreasing  
217 genetic diversity in breeding pool, and clonal propagation. The deleterious mutations  
218 should be important targets for future cassava breeding programs. Genomic selection  
219 and genomic editing technologies<sup>21</sup> are anticipated to help purge deleterious mutations  
220 and improve this globally important crop.



## 221 ONLINE METHODS

### 222 Samples and whole genome sequencing

223 To maximize the diversity and representation for cassava, all samples were selected  
224 based on breeders' choice and diversity analysis from accessions included in Next  
225 Generation Cassava Breeding project ([www.nextgencassava.org](http://www.nextgencassava.org)). Whole genome  
226 sequences were generated from 241 cassava accessions including 203 elite breeding  
227 accessions, 16 progenitors (*M. falbellifolia*, *M. peruviana*)<sup>7</sup>, 11 hybrid/tree cassava  
228 accessions and 11 wild relative cassava accessions (*M. glaziovii* and others)  
229 (**Supplementary Table 1**). Among 241 cassava accessions, 172 accessions were  
230 sequenced at the Genomic Diversity Facility at Cornell University, Ithaca, NY, USA.  
231 Standard Illumina PCR-free libraries were constructed with insert size of 500-bp using  
232 Illumina standard protocol. Sequences of 200-bp length were generated using Illumina  
233 HiSeq 2500 and 150-bp length were generated using NextSeq Series Desktop  
234 sequencers. Donald Danforth Plant Science Center, St. Louis, MO, USA generated ~20x  
235 coverage sequences for 15 elite cassava accessions. Sequences for remaining 54 cassava  
236 accessions were collected from HapMap<sup>10</sup>, generated at the University of California at  
237 Berkeley (USA).

238

### 239 Alignment of reads and variant calling for generation of cassava haplotype map 240 (HapMapII)

241 The cassava genome was found to have large amounts of repeat sequences<sup>10</sup>. To  
242 minimize misalignment, these repeats were pre-filtered by aligning the sequences to a  
243 bait containing repeat sequences and organelle sequences (**Supplementary Fig. 1**).  
244 Remaining sequences after pre-filtering were aligned to reference genome (v6.1) using  
245 burrows-wheeler alignment with maximal exact matches (BWA-MEM) algorithm  
246 (<http://bio-bwa.sourceforge.net/bwa.shtml#13>). To ensure high quality SNP calling,  
247 especially for those rare variants, we developed an in-house pipeline, FastCall  
248 (<https://github.com/Fei-Lu/FastCall>), to perform the stringent variant discovery. The  
249 procedures include: 1) Genomic positions having both insertion and deletion variants  
250 were ignored, since these sites were likely in complex regions with many misalignments;  
251 2) For multiple allelic sites, if the third allele had more than 20% depth in any individual,  
252 the site was ignored; 3) For a specific site, if the minor allele did not have a depth  
253 between 40% and 60% in at least one individual when individual depth was greater than  
254 5, the site was ignored; 4) A chi square test for allele segregation<sup>13</sup> in all individual is  
255 performed. The sites with *P*-value more than  $1.0 \times 10^{-03}$  were ignored. 5) On average,  
256 over 30X depth was used to for individual genotype calls. The genotype likelihood was  
257 calculated based on multinomial test reported by Hohenlohe *et. al*<sup>32</sup>. To remove  
258 potential spurious variants arising from paralogs, an additional filter was applied to keep  
259 only variants with depth between 7,500 and 11,500. The missing data was about 4%.  
260 The genotypes were imputed and phased into haplotypes using BEAGLE v4.1<sup>12</sup>. A total  
261 of 10% of the genotypes were masked before imputation to calculate the imputation  
262 accuracy.

263

### 264 Population genetics analysis

265 SNP density, pair-wise nucleotide diversity ( $\pi$ ), Tajima's D and  $F_{st}$  were calculated using  
266 VCFtools<sup>33</sup> (**Supplementary Fig. 8**). Principal component analysis was carried out in Trait  
267 Analysis by aSSociation, Evolution and Linkage (TASSEL)<sup>34</sup>. Recombination rates were  
268 obtained from cassava HapMap1 source<sup>10</sup>.

269

### 270 **Genomic evolutionary rate profiling (GERP)**

271 Constrained portion of cassava genome was identified by quantifying rejected  
272 substitutions (strength of purifying selection) using GERP++ program<sup>16</sup>. Multiple whole  
273 genome sequence alignment was carried out for the seven species in Malpighiales clade  
274 of plant kingdom, including cassava, rubber (*Hevea brasiliensis*), jatropha (*Jatropha*  
275 *curcas*), castor bean (*Ricinus communis*), willow (*Salix purpurea*), flax (*Linum*  
276 *usitatissimum*), and poplar (*Populus trichocarpa*). Phylogenetic tree and neutral branch  
277 length (estimated from 4-fold degenerate sites) were used to quantify constraint  
278 intensity at every position on cassava genome. Cassava genome sequence was  
279 eliminated during the site specific observed estimates (RS scores) to eliminate the  
280 confounding influence of deleterious derived alleles segregating in cassava populations  
281 that are present in reference sequence.

282

### 283 **Identifying deleterious mutation**

284 Amino acid substitution and their effects on protein function were predicted using  
285 'Sorting Tolerant From Intolerant (SIFT)' algorithm<sup>19</sup>. Non-synonymous mutations with  
286 SIFT score < 0.05 were defined as putative deleterious mutations. SIFT (< 0.05) and GERP  
287 (>2) annotations were combined to identify the deleterious mutations existing in  
288 constrained portion of the genome. These deleterious mutations were used to calculate  
289 mutation burden of cassava.

290

### 291 **Identifying selective sweep regions**

292 Cross-population composite likelihood approach (XP-CLR) method<sup>27</sup> was used to identify  
293 the selective sweeps in two contrasts: Latin America cassava accessions (test  
294 populations) against progenitors (*M. esc. ssp flabellifolia*, reference population) for  
295 domestication event and African cassava accessions (test populations) against Latin  
296 American cassava accessions (reference population) for recent improvement in Africa.  
297 Selection scan was performed across the genome using 0.5 cM sliding window between  
298 the SNPs spacing of 2-kb. A genetic map of cassava generated by International Cassava  
299 Genetic Map Consortium<sup>35</sup> was used in the XP-CLR analysis. XP-CLR scores were  
300 normalized using Z-score and smoothed spline technique with R-package (GenWin)<sup>36</sup>.  
301 Outlier peaks were selected which were above than 99 percentile of normalized values.  
302 AgriGO<sup>37</sup> and REVIGO<sup>38</sup> tools were used for GO enrichment analysis.

303

### 304 **Mutation burden in cassava accessions**

305 Number of derived deleterious alleles present in each cassava accessions were counted  
306 to identify mutation burden in cassava accessions in three models (homozygous  
307 mutation burden, heterozygous mutation burden, and total mutation burden).  
308 Homozygous mutation burden is the number of derived deleterious alleles in the



309 homozygous state. Heterozygous mutation burden is the number of derived deleterious  
310 alleles existing in the heterozygous state. Total mutation burden is the number of  
311 derived deleterious alleles existing in an accession (2 x homozygous mutation burden +  
312 heterozygous mutation burden)<sup>15,39</sup>.

313

#### 314 **Comparison of observed and expected mutation burden under HWE**

315 A bootstrap approach (with replacement) was used to resample cultivated cassava  
316 accessions from both Latin American (24 samples) and African (174 samples) breeding  
317 pools. The process was repeated for 10,000 times to generate the distribution of  
318 expected homozygous and heterozygous mutation burden. For each resampling,  
319

319

$$320 \quad b_{ho} = \sum_{i=1}^n d_i^2, \quad b_{he} = \sum_{i=1}^n 2(1 - d_i)d_i$$

321

322 where  $b_{ho}$  is the expected homozygous mutation burden under HWE,  $b_{he}$  is the expected  
323 heterozygous mutation burden under HWE,  $n$  is the total number of deleterious  
324 mutations identified ( $n = 22,495$ ),  $d_i$  is the allele frequency of  $i$ th deleterious allele in the  
325 sampled population. The observed mutation burden was calculated for each accession  
326 as described in 'mutation burden in cassava accessions'. The means of observed  
327 homozygous and heterozygous mutation were used for the comparison.

328

329

330 Data access:

331 Whole genome sequences, raw and imputed HapMapII SNPs can be accessed from  
332 CassavaBase at <ftp://ftp.cassavabase.org/HapMapII/>.

333

334

#### 335 **ACKNOWLEDGEMENTS**

336 This work was supported by the Bill & Melinda Gates Foundation (BMGF:  
337 #01511000147), with additional support from NSF Plant Genome Research Project  
338 (#1238014) and the USDA-ARS. We thank Next Generation cassava project  
339 ([www.nextgencassava.org](http://www.nextgencassava.org)) for helping us to choose the accessions to include in whole  
340 genome sequencing efforts. We thank Simon E. Prochnik (DOE Joint Genome Institute,  
341 Walnut Creek, CA, USA) for his timely help during the analysis.

342

#### 343 **AUTHORS CONTRIBUTIONS**

344 The manuscript was prepared by P.R., F.L.. Data analysis was carried out by P.R., F.L. and  
345 E.S.B.. Whole genome sequences for 54 accessions included in HapMap<sup>10</sup> are provided  
346 by J.V.B. W.E., I.Y.R., C.E., R.K. and R.S.B. provided the germplasm for WGS. All authors  
347 provided their comments and edited the manuscript. F.L. and E.S.B designed and  
348 coordinated the project.

349

#### 350 **COMPETING FINANCIAL INTERESTS**

351 The authors declare no competing financial interests.

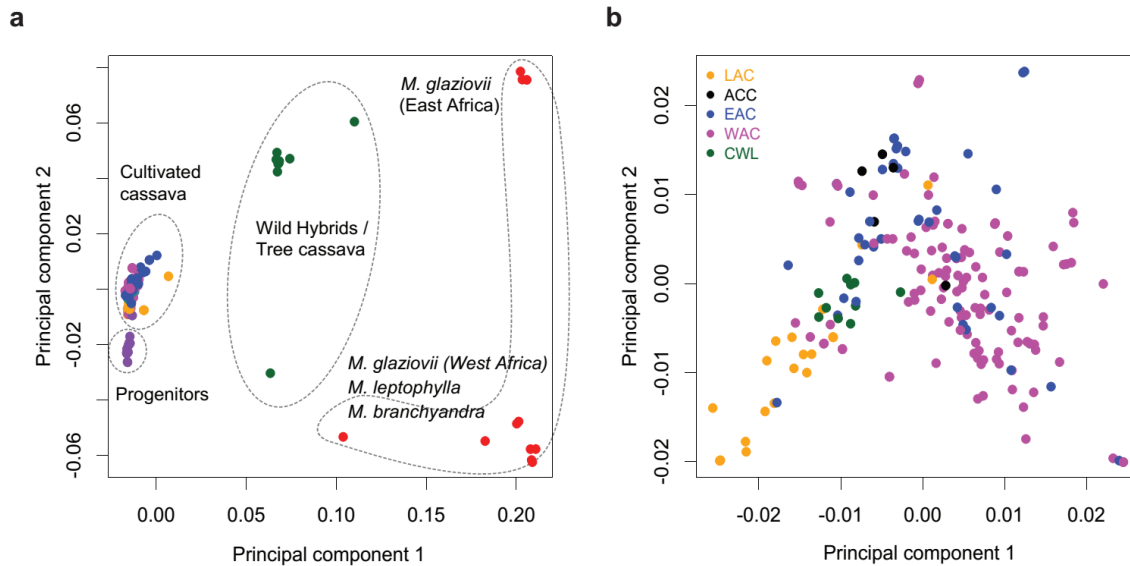
## 352 References

- 353 1. Raven, P., Fauquet, C., Swaminathan, M.S., Borlaug, N. & Samper, C. Where Next  
354 for Genome Sequencing? *Science* 311, 468-468 (2006).
- 355 2. Olsen, K.M. & Schaal, B.A. Evidence on the origin of cassava: Phylogeography of  
356 *Manihot esculenta*. *Proceedings of the National Academy of Sciences* 96, 5586-  
357 5591 (1999).
- 358 3. Wang, W. et al. Cassava genome from a wild ancestor to cultivated varieties. *Nat*  
359 *Commun* 5(2014).
- 360 4. McDonald, M.J., Rice, D.P. & Desai, M.M. Sex speeds adaptation by altering the  
361 dynamics of molecular evolution. *Nature* 531, 233-236 (2016).
- 362 5. McKey, D., Elias, M., Pujol, B. & Duputié, A. The evolutionary ecology of clonally  
363 propagated domesticated plants. *New Phytologist* 186, 318-332 (2010).
- 364 6. Charlesworth, D. & Willis, J.H. The genetics of inbreeding depression. *Nat Rev*  
365 *Genet* 10, 783-796 (2009).
- 366 7. Rojas, M.C. et al. Analysis of Inbreeding Depression in Eight S1 Cassava Families.  
367 *Crop Science* 49, 543-548 (2009).
- 368 8. Nuwamanya, E., Herselman, L. & Ferguson, M. Segregation of selected  
369 agronomic traits in six S1 cassava families. *Journal of Plant Breeding and Crop*  
370 *Science* 3, 154-160 (2011).
- 371 9. Allem, A.C. The closest wild relatives of cassava ( *Manihot esculenta* Crantz).  
372 *Euphytica* 107, 123-133 (1999).
- 373 10. Bredeson, J.V. et al. Sequencing wild and cultivated cassava and related species  
374 reveals extensive interspecific hybridization and genetic diversity. *Nat Biotech*  
375 34, 562-570 (2016).
- 376 11. Prochnik, S. et al. The Cassava Genome: Current Progress, Future Directions.  
377 *Tropical Plant Biology* 5, 88-94 (2012).
- 378 12. Browning, Brian L. & Browning, Sharon R. Genotype Imputation with Millions of  
379 Reference Samples. *The American Journal of Human Genetics* 98, 116-126.
- 380 13. Chia, J.-M. et al. Maize HapMap2 identifies extant variation from a genome in  
381 flux. *Nat Genet* 44, 803-807 (2012).
- 382 14. Tennessen, J.A. et al. Evolution and Functional Impact of Rare Coding Variation  
383 from Deep Sequencing of Human Exomes. *Science* 337, 64-69 (2012).
- 384 15. Fu, W. et al. Analysis of 6,515 exomes reveals the recent origin of most human  
385 protein-coding variants. *Nature* 493, 216-220 (2013).
- 386 16. Davydov, E.V. et al. Identifying a High Fraction of the Human Genome to be  
387 under Selective Constraint Using GERP++. *PLoS Comput Biol* 6, e1001025 (2010).
- 388 17. Rodgers-Melnick, E. et al. Recombination in diverse maize is stable, predictable,  
389 and associated with genetic load. *Proceedings of the National Academy of*  
390 *Sciences* 112, 3823-3828 (2015).
- 391 18. Mackay, T.F.C. et al. The *Drosophila melanogaster* Genetic Reference Panel.  
392 *Nature* 482, 173-178 (2012).
- 393 19. Kumar, P., Henikoff, S. & Ng, P.C. Predicting the effects of coding non-  
394 synonymous variants on protein function using the SIFT algorithm. *Nat. Protocols*  
395 4, 1073-1081 (2009).

- 396 20. Mezmouk, S. & Ross-Ibarra, J. The Pattern and Distribution of Deleterious  
397 Mutations in Maize. *G3: Genes/Genomes/Genetics* 4, 163-171 (2014).
- 398 21. Horvath, P. & Barrangou, R. CRISPR/Cas, the Immune System of Bacteria and  
399 Archaea. *Science* 327, 167-170 (2010).
- 400 22. Keller, P.J. & Knop, M. Evolution of Mutational Robustness in the Yeast Genome:  
401 A Link to Essential Genes and Meiotic Recombination Hotspots. *PLoS Genet* 5,  
402 e1000533 (2009).
- 403 23. Hussin, J.G. et al. Recombination affects accumulation of damaging and disease-  
404 associated mutations in human populations. *Nat Genet* 47, 400-404 (2015).
- 405 24. Haddrill, P.R., Halligan, D.L., Tomaras, D. & Charlesworth, B. Reduced efficacy of  
406 selection in regions of the Drosophila genome that lack crossing over. *Genome*  
407 *Biology* 8, 1-9 (2007).
- 408 25. Jørgensen, K. et al. Cassava Plants with a Depleted Cyanogenic Glucoside  
409 Content in Leaves and Tubers. Distribution of Cyanogenic Glucosides, Their Site  
410 of Synthesis and Transport, and Blockage of the Biosynthesis by RNA  
411 Interference Technology. *Plant Physiology* 139, 363-374 (2005).
- 412 26. Conn, E.E. Cyanogenic Compounds. *Annual Review of Plant Physiology* 31, 433-  
413 451 (1980).
- 414 27. Chen, H., Patterson, N. & Reich, D. Population differentiation as a test for  
415 selective sweeps. *Genome Research* 20, 393-402 (2010).
- 416 28. Ros, R., Muñoz-Bertomeu, J. & Krueger, S. Serine in plants: biosynthesis,  
417 metabolism, and functions. *Trends in Plant Science* 19, 564-569 (2014).
- 418 29. Benstein, R.M. et al. Arabidopsis Phosphoglycerate Dehydrogenase1 of the  
419 Phosphoserine Pathway Is Essential for Development and Required for  
420 Ammonium Assimilation and Tryptophan Biosynthesis. *The Plant Cell* 25, 5011-  
421 5029 (2013).
- 422 30. Marsden, C.D. et al. Bottlenecks and selective sweeps during domestication have  
423 increased deleterious genetic variation in dogs. *Proceedings of the National*  
424 *Academy of Sciences* 113, 152-157 (2016).
- 425 31. Agrawal, A.F. & Whitlock, M.C. Inferences About the Distribution of Dominance  
426 Drawn From Yeast Gene Knockout Data. *Genetics* 187, 553-566 (2011).
- 427 32. Hohenlohe, P.A. et al. Population Genomics of Parallel Adaptation in Threespine  
428 Stickleback using Sequenced RAD Tags. *PLoS Genet* 6, e1000862 (2010).
- 429 33. Danecek, P. et al. The variant call format and VCFtools. *Bioinformatics* 27, 2156-  
430 2158 (2011).
- 431 34. Bradbury, P.J. et al. TASSEL: software for association mapping of complex traits  
432 in diverse samples. *Bioinformatics* 23, 2633-2635 (2007).
- 433 35. Consortium, I.C.G.M. High-Resolution Linkage Map and Chromosome-Scale  
434 Genome Assembly for Cassava (*Manihot esculenta* Crantz) from 10 Populations.  
435 *G3: Genes/Genomes/Genetics* 5, 133-144 (2015).
- 436 36. Beissinger, T.M., Rosa, G.J., Kaeppler, S.M., Gianola, D. & de Leon, N. Defining  
437 window-boundaries for genomic analyses using smoothing spline techniques.  
438 *Genetics Selection Evolution* 47, 1-9 (2015).

- 439 37. Du, Z., Zhou, X., Ling, Y., Zhang, Z. & Su, Z. agriGO: a GO analysis toolkit for the  
440 agricultural community. *Nucleic Acids Research* 38, W64-W70 (2010).
- 441 38. Supek, F., Bošnjak, M., Škunca, N. & Šmuc, T. REVIGO Summarizes and Visualizes  
442 Long Lists of Gene Ontology Terms. *PLoS ONE* 6, e21800 (2011).
- 443 39. Henn, B.M. et al. Distance from sub-Saharan Africa predicts mutational load in  
444 diverse human genomes. *Proceedings of the National Academy of Sciences* 113,  
445 E440-E449 (2016).
- 446
- 447

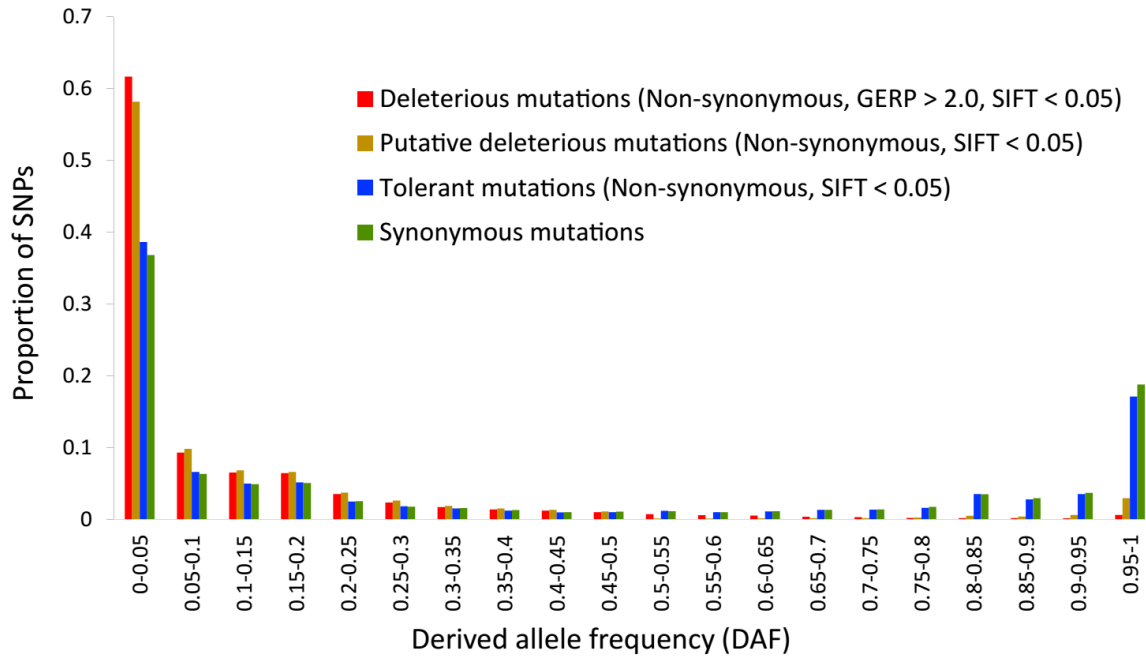
448 **Figures**



449

450

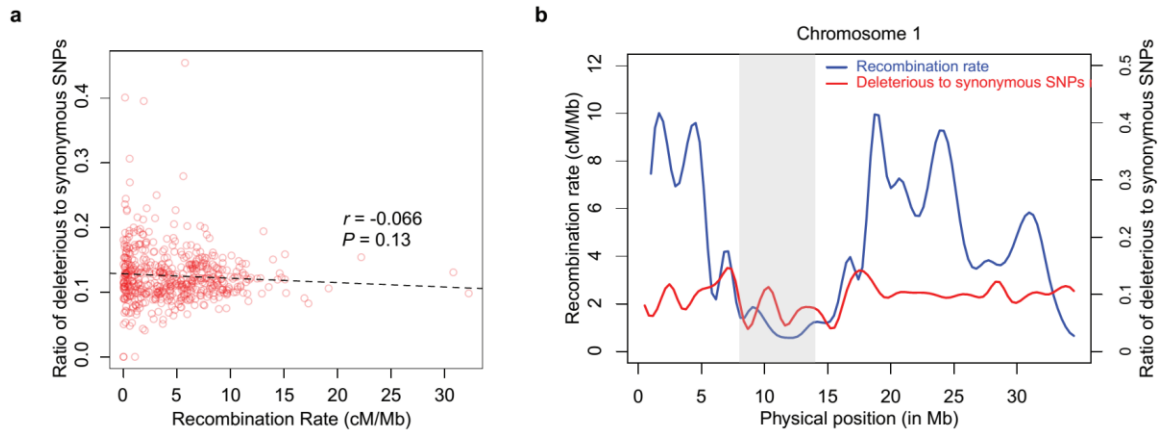
451 **Figure 1** Principal component analysis (PCA) of cassava accessions included in cassava  
452 HapMapII. (a) PCA of all cassava accessions (progenitors, cultivated, and wild cassava  
453 accessions). A total of 43.8% genetic variance is captured in first two principal  
454 components. (b) PCA of cultivated cassava clones. A total of 9.1% genetic variance is  
455 captured in first two principal components. The abbreviations are represented as  
456 follows: LAC – Latin American cassava, ACC – Asian Cultivated cassava, EAC – East  
457 African cassava, WAC – West African cassava, CWL – Crosses between WAC and LAC.  
458



459  
460  
461  
462  
463  
464  
465

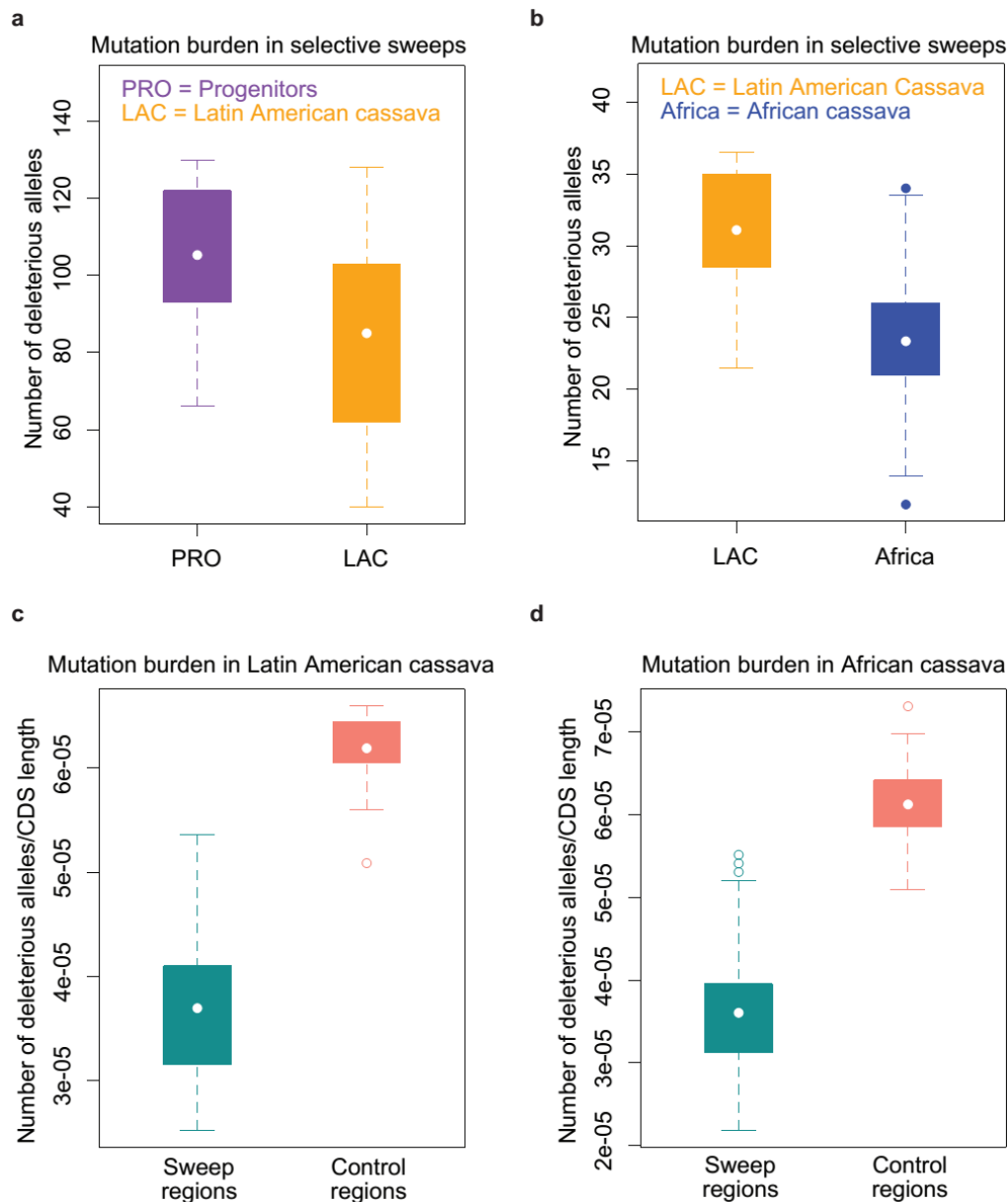
**Figure 2** Site allele frequency spectrum of deleterious mutations in cassava genome. Derived allele frequency (DAF) distribution of alleles are presented. Rubber genome is used as the out group to define derived alleles.





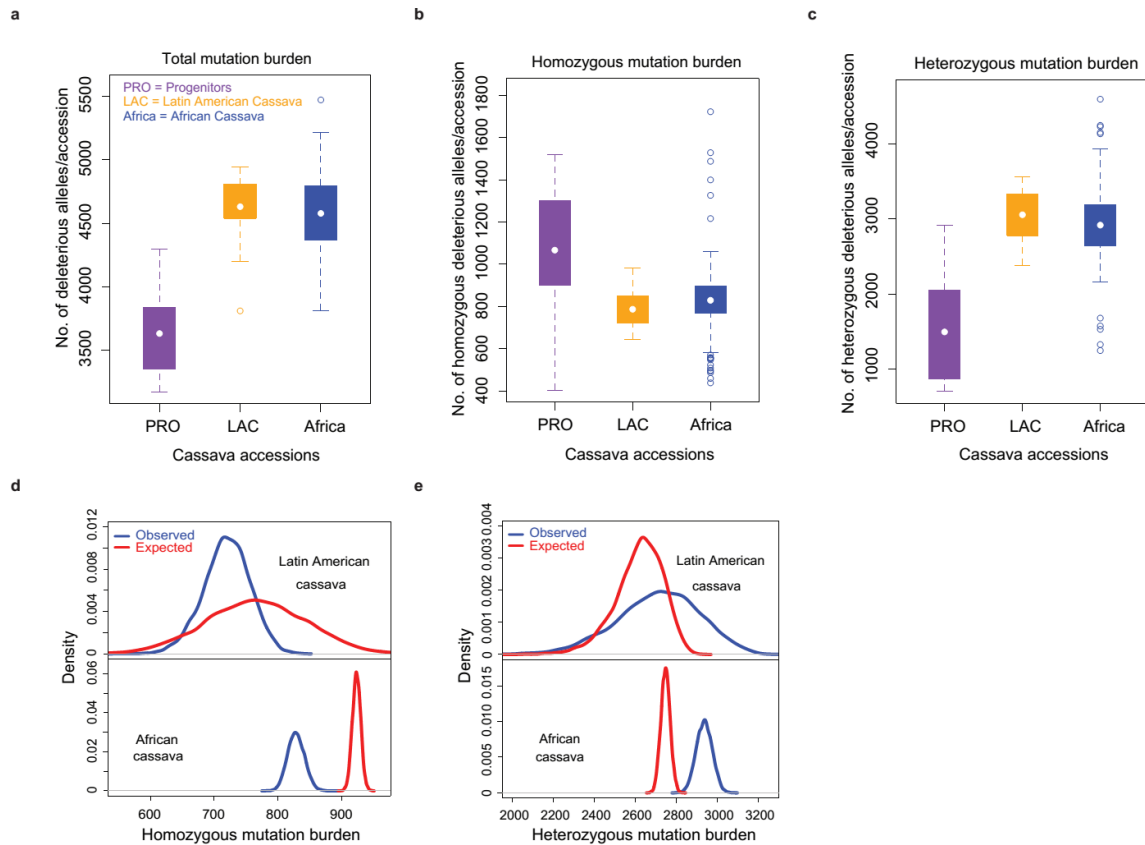
466  
467  
468  
469  
470  
471  
472

**Figure 3** Effect of recombination on the distribution of deleterious mutations in cassava genome. (a) Correlation between recombination rate and number of deleterious mutations in the genome. (b) Distribution of deleterious mutations as a function of recombination rate on chromosome 1.



473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483

**Figure 4** Mutation burden in selective sweep regions. (a) Mutation burden between progenitors and Latin American cassava accessions in domestication sweep regions. (b) Mutation burden between Africa and Latin American cassava accessions in sweep regions identified in recent improvement in Africa. (c) Mutation burden in Latin American cassava accessions between domestication selective sweeps and control regions (rest of the genome). (d) Mutation burden in African cassava accessions between sweep regions identified in recent improvement and control regions (rest of the genome) in Africa.



484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497

**Figure 5** Mutation burden in cassava populations. (a) Total mutation burden in progenitors, Latin American cassava and African cassava accessions. Bottleneck during domestication increased mutation burden. Demography in Africa has no significant influence on mutation burden in African cassava accessions. (b) Homozygous mutation burden in cassava populations. Domestication decreased homozygous mutation burden in cultivated cassava. (c) Heterozygous mutation burden in cassava populations. Domestication increased heterozygous mutation burden in cultivated cassava. (d) Comparison between the observed homozygous mutation burden and the expected homozygous mutation burden under HWE assumption in cultivated cassava. (e) Comparison between the observed heterozygous mutation burden and the expected heterozygous mutation burden under HWE assumption in cultivated cassava.