

1 **Cassava HapMap: Managing genetic load 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>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 30 million segregating variants. While domestication**  
22 **modified starch and ketone metabolism pathways for human consumption, the**  
23 **concomitant bottleneck and clonal propagation resulted in a large proportion of fixed**  
24 **deleterious amino acid changes, raised the number of deleterious mutations by 24%,**  
25 **and shifted the mutational burden towards common variants. Deleterious mutations**  
26 **are ineffectively purged due to limited recombination in cassava genome. Recent**  
27 **breeding efforts maintained the yield by masking the harmful effects of deleterious**  
28 **mutations through shielding the most damaging recessive mutations in the**  
29 **heterozygous state, but unable to purge the load, which should be a key target for**  
30 **future cassava breeding.**

31

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

44 in elite cassava accessions, where a single generation of inbreeding results in >60%  
45 reduction in fresh root yield<sup>7,8</sup>. In this study, we aimed to identify deleterious mutations  
46 in cassava populations, which in turn can help accelerate cassava breeding by allowing  
47 breeders to purge deleterious mutations more efficiently.

48 We conducted a comprehensive characterization of genetic variation by whole genome  
49 sequencing (WGS) of 241 cassava accessions, including 203 elite breeding accessions (*M.*  
50 *esculenta* Crantz), 16 close relatives (*M. esculenta* ssp. *flabellifolia*, *M. esculenta* ssp.  
51 *peruviana*) of modern cultivars<sup>2,9</sup>, 11 hybrid/tree cassava accessions, and 11 more  
52 divergent wild relatives (*M. glaziovii* and others) (**Supplementary Table 1**). Samples  
53 included 54 accessions from an initial haplotype map I (HapMapI) study<sup>10</sup>. Wild *M.*  
54 *glaziovii* has been used extensively in cassava breeding programs to transfer disease  
55 resistance alleles to cultivated cassava (e.g., Amani Breeding program)<sup>8</sup>. On average,  
56 more than 30x coverage sequences were generated for each accession. The 518.5 Mb  
57 cassava genome (V6.1) has roughly 51% repetitive elements with several common  
58 recent retrotransposons<sup>10</sup>. To exclude misalignment and ensure high quality of variant  
59 calling, repeat sequences were pre-filtered using repeat bait (**Supplementary Fig. 1**) and  
60 the remaining sequences were aligned against the cassava reference genome v6.1<sup>10,11</sup>.  
61 Variants from low copy regions of the genome were identified to develop the cassava  
62 haplotype map II (HapMapII) with 30.5 million variants (28.38 million SNPs and 2.14  
63 million indels) and with a low error rate of 0.02%, which is the proportion of segregating  
64 sites in the reference accession (**Supplementary Fig. 2**). Cultivated cassava exhibited  
65 12.18 million variants (**Supplementary Table 2**), of which more than 50% were found to  
66 be rare (<5% minor allele frequency (MAF)) (**Supplementary Table 2 and**  
67 **Supplementary Fig. 3**). Haplotypes were phased and missing genotypes were imputed  
68 with high accuracy using BEAGLE v4.1<sup>12</sup> (accuracy  $r^2 = 0.989$ ) (**Supplementary Fig. 4**).  
69 Linkage disequilibrium was as low as in maize<sup>13</sup> and decayed to an average  $r^2 = 0.1$  in  
70 3,000 bp (**Supplementary Fig. 5**).

71 Cultivated cassava presented lower nucleotide diversity ( $\pi=0.0041$ ) compared with its  
72 progenitors (*M. esc.* ssp. *flabellifolia*,  $\pi=0.0057$ ). In addition, a close relationship  
73 between the two species was observed from phylogenetic analysis (**Supplementary Fig.**  
74 **6**). 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 ( $F_{st}$ : 0.15), and high  
79 genetic differentiation from tree cassava ( $F_{st}$ : 0.31) and wild relatives ( $F_{st}$ : 0.43)  
80 (**Supplementary Table 2 and Supplementary Figs. 7 and 8**). 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-Mbp of the  
92 genome (20%) of cassava was constrained (GERP score > 0) (**Supplementary Fig. 9**). The  
93 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. 10**). 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 2.7% of coding SNPs in cultivated cassava were non-synonymous mutations  
100 (**Supplementary Table 2** and **Supplementary Fig. 11**), of which 20% (68,894) 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 23,697 deleterious mutations (**Supplementary Fig. 12**).

104 To estimate the individual mutation load, 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 (6%, **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 153  
110 fixed deleterious mutations. These deleterious mutations are not targets for standard  
111 breeding as they do not segregate, but they are the potential targets for genome  
112 editing<sup>21</sup>. Together with the other 23,544 segregating deleterious mutations, the  
113 mutational load in cassava was substantial. Given the several millennia of breeding in  
114 the species, why are these deleterious mutations still in cultivated cassava and how  
115 were breeders managing them? We evaluated the effects of recombination, selection,  
116 and drift, as the main processes controlling the distribution of deleterious mutations in  
117 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 ( $r=-0.07$ ,  $P =$   
122  $0.081$ , **Fig. 3a**). Deleterious mutation were nearly uniformly spread across the cassava  
123 genome (**Fig. 3b** and **Supplementary Fig. 13**), 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 genetic load 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 224 selective  
137 sweeps containing 484 genes in Latin American accessions (**Supplementary Fig. 14**).  
138 Genes in these sweep regions were enriched for starch and sucrose synthesis (3.6-fold  
139 enrichment; FDR =  $4.1 \times 10^{-04}$ ) and cellular ketone metabolism (3.3-fold enrichment; FDR  
140 =  $8.1 \times 10^{-05}$ ) (**Supplementary Fig. 14**). 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 286 selective sweeps were identified  
144 containing 470 genes. These genes were enriched for amino acid metabolism (5.8-fold  
145 enrichment, FDR =  $5.5 \times 10^{-06}$ ) and stimulus response (3.4-fold enrichment, FDR =  $9.2 \times$   
146  $10^{-04}$ , **Supplementary Fig. 15**), reflecting that disease resistance accessions were  
147 selected in recent breeding program in Africa<sup>8</sup>.

148 How was the genetic load shaped in the selective sweeps? We found that Latin  
149 American accessions showed 21% less ( $P = 0.006$ , **Fig. 4a**) deleterious mutations than  
150 progenitors in sweep regions. Similarly, African accessions exhibited a 26% drop ( $P = 1.8$   
151  $\times 10^{-06}$ , **Fig. 4b**) in sweeps compared to Latin American accessions. In addition to the  
152 comparison between populations, significant reductions of deleterious mutations were  
153 observed within population by comparing sweep regions and the rest of the genome.  
154 For example, selective sweeps presented 44% depletion ( $P = 7 \times 10^{-13}$ , **Fig. 4c**) of  
155 deleterious mutations in Latin American accessions and 23% reduction ( $P = 4 \times 10^{-63}$ , **Fig.**  
156 **4d**) in African accessions. This implies that haplotypes containing fewer deleterious  
157 mutations were favored during selection.

158 However, drift after domestication played a more important role in affecting mutational  
159 load in cassava. Although Latin American accessions and African accessions had a similar  
160 number of deleterious mutations ( $P = 0.76$ , **Fig. 5a**), they presented a prominent  
161 increase of total load by 24% ( $P = 9.4 \times 10^{-09}$ , **Fig. 5a**) when compared with progenitors,  
162 and shifted the mutational burden towards common deleterious variants (**Fig. 5b**). The  
163 increase of deleterious mutations during domestication was also found in dog<sup>28</sup>. The  
164 results suggest that the severe bottleneck of domestication and shift from sexual

165 reproduction to clonal propagation resulted in a rapid accumulation of deleterious  
166 mutations in cultivated cassava.

167 How have the breeders been able to maintain yield, given the substantial growth of  
168 mutational load in cultivated cassava? This became apparent when the homozygous  
169 deleterious mutations and heterozygous deleterious mutations were compared.  
170 Relative to *M. esculenta* ssp. *flabellifolia*, the homozygous genetic load substantially  
171 decreased by 23% ( $P = 6 \times 10^{-03}$ , **Fig. 5c**) in cultivated accessions, while the heterozygous  
172 load remarkably increased by 82% ( $P = 8 \times 10^{-07}$ , **Fig. 5d**), despite the reduced genetic  
173 diversity in cultivated cassava. This suggests that breeders have been trying to manage  
174 the recessive deleterious mutations in the heterozygous state to mask their harmful  
175 effects. Mutations with large homozygous effect are more likely to be recessive<sup>29</sup>. We  
176 found nearly 61% of deleterious mutations occurred only in the heterozygous state (**Fig.**  
177 **5e**). These were likely to be the lethal/strong deleterious mutations, resulting in the  
178 significant yield loss in the first generation of selfed cassava plants<sup>7,8</sup>.

179 Cassava is a major staple crop feeding hundreds of millions people. Using deep  
180 sequencing of a comprehensive and representative collection of 241 cassava accessions,  
181 we developed the HapMapII, a highly valuable resource for cassava genetic studies and  
182 breeding. In this vegetatively propagated species, deleterious mutations have been  
183 accumulating rapidly due to the lack of recombination. The bottleneck event during  
184 domestication exacerbated the existing genetic load in cassava. Breeding efforts  
185 successfully maintained the yield by selecting high fitness haplotypes at a few hundred  
186 loci and handling most damaging mutations in the heterozygous state. However,  
187 breeders were unable to purge the load due to limited recombination, instead they  
188 shielded deleterious mutations by increasing the heterozygosity while screening  
189 thousands of potential hybrids (**Supplementary Fig. 16**). In the short term, this practice  
190 for managing genetic load may produce gains in yield. In the long run, however, a  
191 mutational meltdown may be triggered by new mutations, decreasing genetic diversity  
192 in breeding pool, and clonal propagation. The deleterious mutations should be  
193 important targets for future cassava breeding programs. Genomic selection and  
194 genomic editing technologies<sup>21</sup> are anticipated to help purge deleterious mutations and  
195 improve this globally important crop.



## 196 ONLINE METHODS

### 197 Samples and whole genome sequencing

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

212

### 213 Alignment of reads and variant calling for generation of cassava haplotype map 214 (HapMapII)

215 The cassava genome was found to have large amounts of repeat sequences<sup>10</sup>. To  
216 minimize misalignment, these repeats were pre-filtered by aligning the sequences to a  
217 bait containing repeat sequences and organelle sequences (**Supplementary Fig. 1**).  
218 Remaining sequences after pre-filtering were aligned to reference genome (V6.1) using  
219 burrows-wheeler alignment with maximal exact matches (BWA-MEM) algorithm  
220 (<http://bio-bwa.sourceforge.net/bwa.shtml#13>). To ensure high quality SNP calling,  
221 especially for those rare variants, we developed an in-house pipeline, FastCall  
222 (<https://github.com/Fei-Lu/FastCall>), to perform the stringent variant discovery. The  
223 procedures include: 1) Genomic positions having both insertion and deletion variants  
224 were ignored, since these sites were likely in complex regions with many misalignments;  
225 2) For multiple allelic sites, if the third allele had more than 20% depth in any individual,  
226 the site was ignored; 3) For a specific site, if the minor allele did not have a depth  
227 between 40% and 60% in at least one individual when individual depth was greater than  
228 5, the site was ignored; 4) A chi square test for allele segregation<sup>13</sup> in all individual is  
229 performed. The sites with *P*-value more than  $1.0 \times 10^{-03}$  were ignored. 5) On average,  
230 over 30X depth was used to for individual genotype calls. The genotype likelihood was  
231 calculated based on multinomial test reported by Hohenlohe *et. al*<sup>30</sup>. The missing data  
232 was about 4%. The genotypes were imputed and phased into haplotypes using BEAGLE  
233 v4.1<sup>12</sup>. A total of 10% of the genotypes were masked before imputation to calculate the  
234 imputation accuracy.

235

### 236 Population genetics analysis

237 SNP density, pair-wise nucleotide diversity ( $\pi$ ), Tajima's *D* and  $F_{st}$  were calculated using  
238 VCFtools<sup>31</sup> (**Supplementary Fig. 8**). Principal component analysis was carried out in Trait

239 Analysis by aSSociation, Evolution and Linkage (TASSEL)<sup>32</sup>. Recombination rates were  
240 obtained from cassava HapMap1 source<sup>10</sup>.

241

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

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

254

### 255 **Identifying deleterious mutation**

256 Amino acid substitution and their effects on protein function were predicted using  
257 'Sorting Tolerant From Intolerant (SIFT)' algorithm<sup>19</sup>. Non-synonymous mutations with  
258 SIFT score < 0.05 were defined as putative deleterious mutations. SIFT (< 0.05) and GERP  
259 (>2) annotations were combined to identify the deleterious mutations existing in  
260 constrained portion of the genome. These deleterious mutations were used to calculate  
261 genetic load of cassava.

262

### 263 **Identifying selective sweep regions**

264 Cross-population composite likelihood approach (XP-CLR) method<sup>27</sup> was used to identify  
265 the selective sweeps in two contrasts: Latin America cassava accessions (test  
266 populations) against progenitors (*M. esc. ssp flabellifolia*, reference population) for  
267 domestication event and African cassava accessions (test populations) against Latin  
268 American cassava accessions (reference population) for recent improvement in Africa.  
269 Selection scan was performed across the genome using 0.5 cM sliding window between  
270 the SNPs spacing of 2-kb. XP-CLR scores were normalized using Z-score and smoothed  
271 spline technique with R-package (GenWin)<sup>33</sup>. Outlier peaks were selected which were  
272 above than 99 percentile of normalized values. AgriGO<sup>34</sup> and REVIGO<sup>35</sup> tools were used  
273 for GO enrichment analysis.

274

### 275 **Genetic load in cassava accessions**

276 Number of derived deleterious alleles present in each cassava accessions were counted  
277 to identify the genetic load in cassava accessions in three models (homozygous load,  
278 heterozygous load, and total load). Homozygous load is the number of derived  
279 deleterious alleles in homozygous state. Heterozygous load is the number of derived  
280 deleterious alleles existing in heterozygous state. Total load is the number of derived  
281 deleterious alleles existing in an accession (2 x homozygous load + heterozygous  
282 load)<sup>15,36</sup>.

283 Data access:

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

286

287

#### 288 ACKNOWLEDGEMENTS

289 This work was supported by the Bill & Melinda Gates Foundation (BMGF:  
290 #01511000147), with additional support from NSF Plant Genome Research Project  
291 (#1238014) and the UDSA-ARS. We thank Next Generation cassava project  
292 ([www.nextgencassava.org](http://www.nextgencassava.org)) for helping us to choose the accessions to include in whole  
293 genome sequencing efforts. We thank Simon E. Prochnik (DOE Joint Genome Institute,  
294 Walnut Creek, CA, USA) for his timely help during the analysis.

295

#### 296 AUTHORS CONTRIBUTIONS

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

302

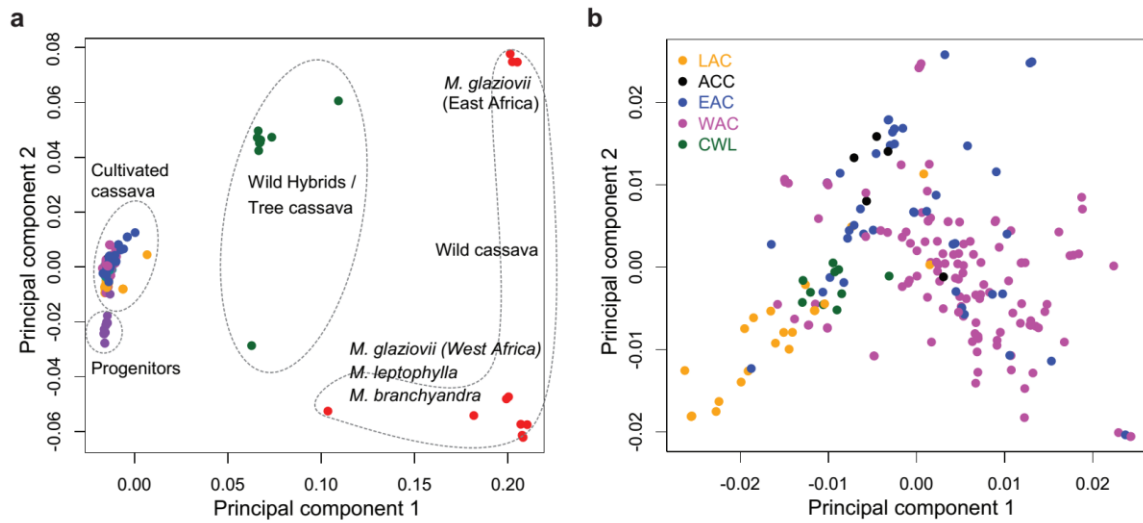
#### 303 COMPETING FINANCIAL INTERESTS

304 The authors declare no competing financial interests.

305



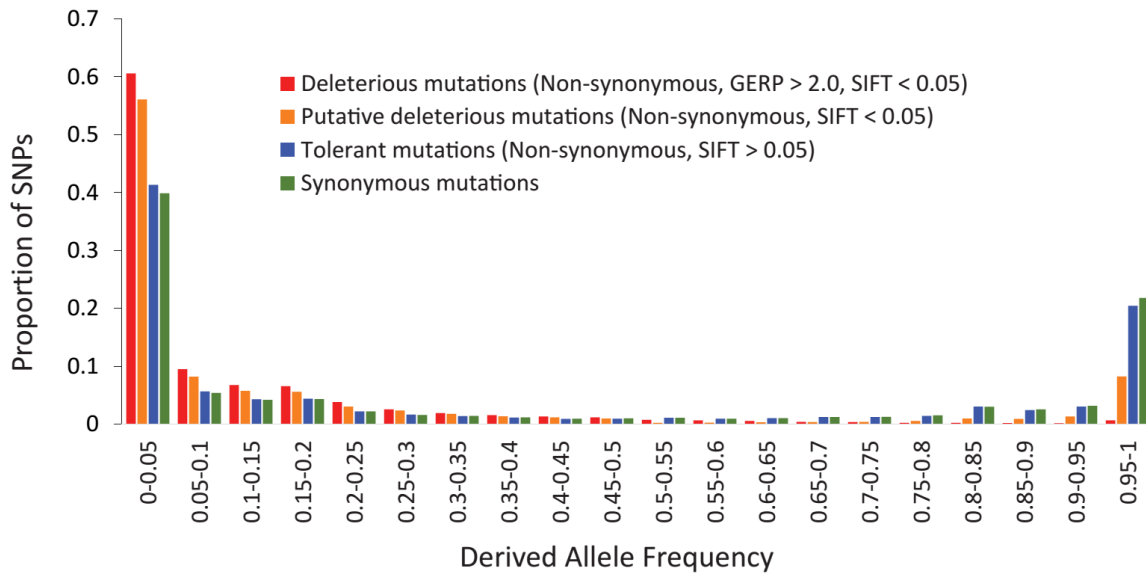
306 **Figures**  
307



308  
309

310 **Figure 1** Principal component analysis (PCA) of cassava accessions included in cassava  
311 HapMapII. (a) PCA of all cassava accessions (progenitors, cultivated, and wild cassava  
312 accessions). A total of 45% genetic variance was captured in first two principal  
313 components. (b) PCA of cultivated cassava clones. A total of 9% genetic variance was  
314 captured in first two principal components. The abbreviations are represented as  
315 follows: LAC – Latin American cassava, ACC – Asian Cultivated cassava, EAC – East  
316 African cassava, WAC – West African cassava, CWL – Crosses between WAC and LAC.  
317

318



319

320

321 **Figure 2** Site allele frequency spectrum of deleterious mutations in cassava genome.

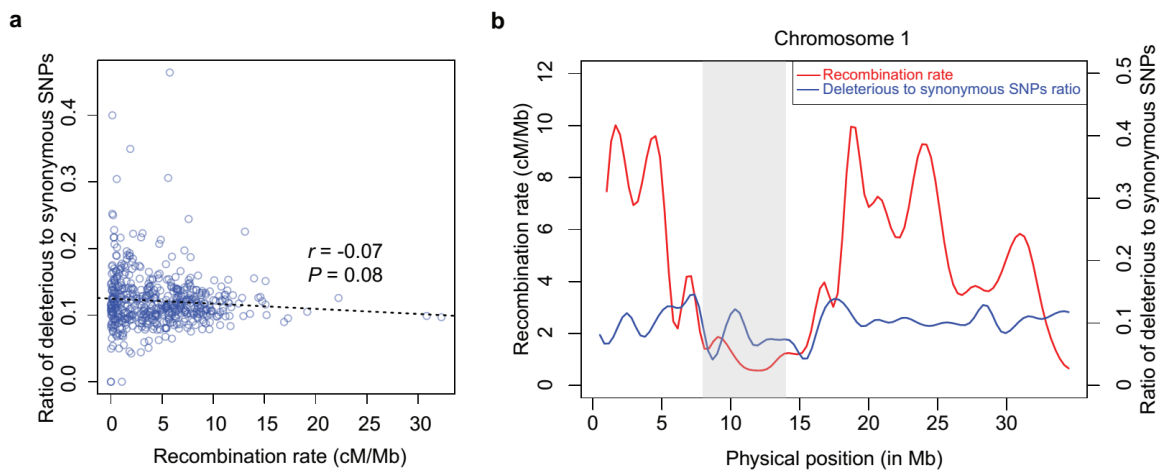
322 Derived allele frequency (DAF) distribution of alleles are presented. Rubber genome is

323 used as the out group to define derived alleles.

324

325

326



327

328

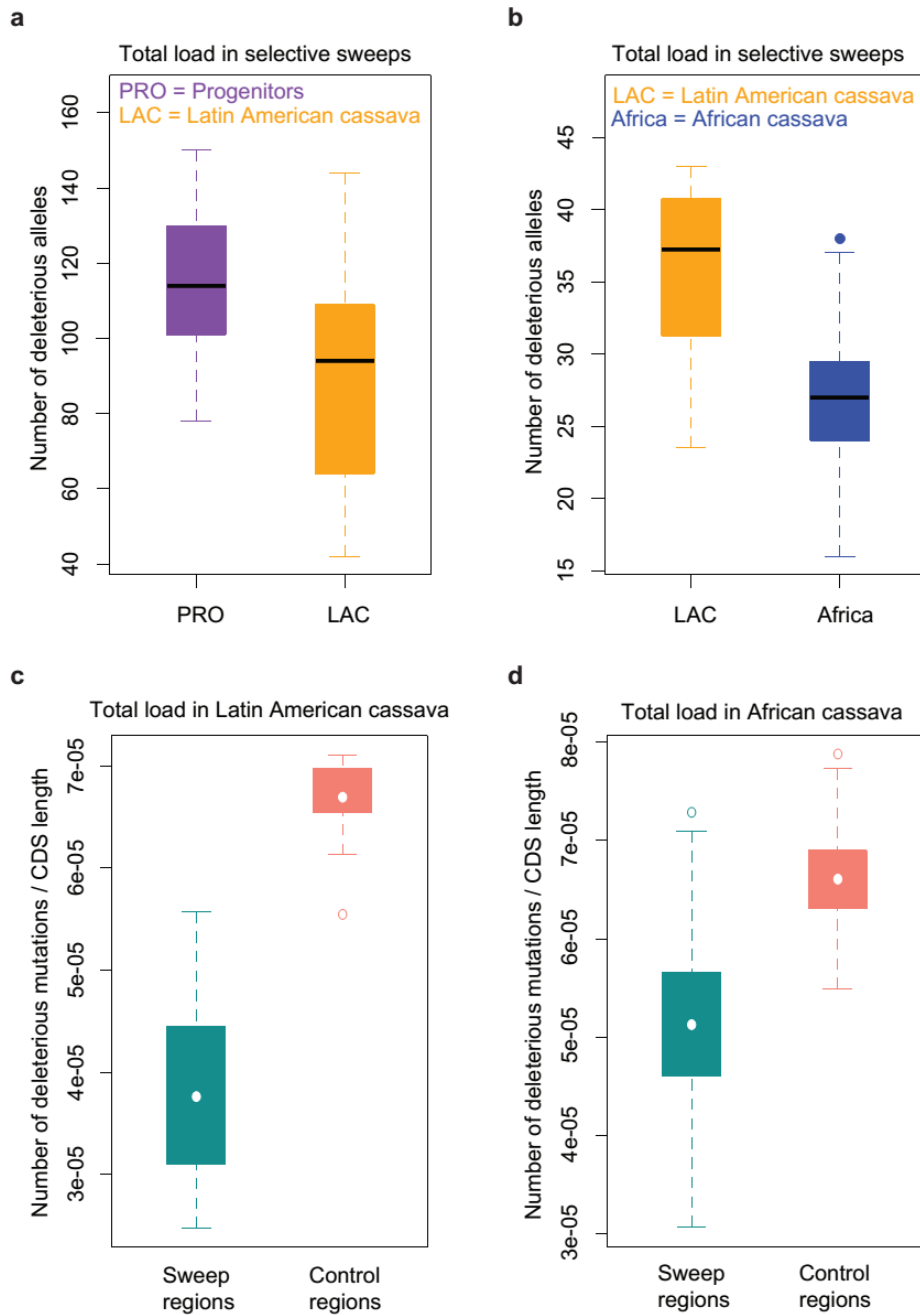
329 **Figure 3** Effect of recombination on the distribution of deleterious mutations in cassava

330 genome. (a) Correlation between recombination rate and number of deleterious

331 mutations in the genome. (b) Distribution of deleterious mutations as a function of

332 recombination rate on chromosome 1.

333



334

335

336

337

338

339

340

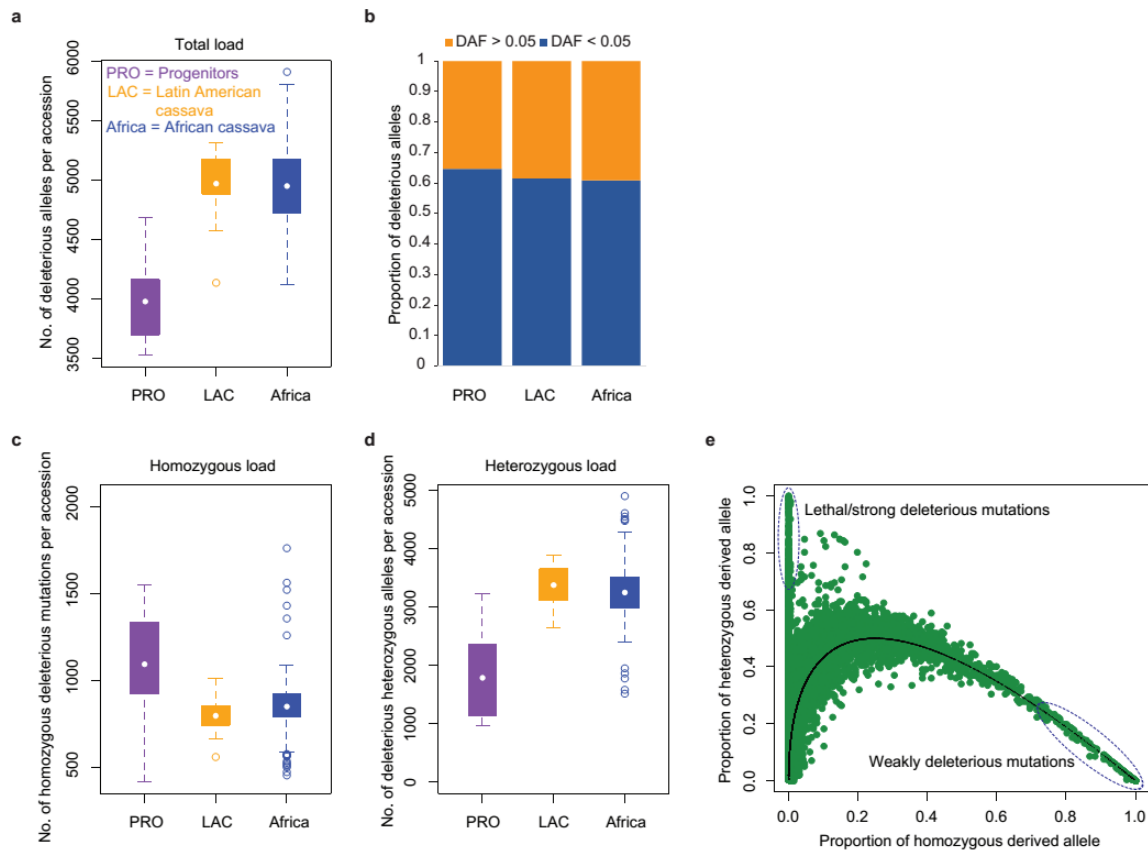
341

342

343

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

344



345

346

347 **Figure 5** Genetic load in cassava populations. (a) Total genetic load in progenitors, Latin  
348 American cassava and African cassava accessions. Bottleneck during domestication  
349 increased the load. Demography in Africa has no significant influence on genetic load in  
350 African cassava accessions. (b) Proportion of deleterious alleles in cassava populations.  
351 (c) Homozygous load in cassava populations. Domestication decreased the homozygous  
352 load in cultivated cassava. (d) Heterozygous load in cassava populations. Domestication  
353 increased the heterozygous load in cultivated cassava. (e) Homozygous and  
354 heterozygous derived allele frequency for deleterious mutations in cultivated cassava  
355 accessions. Black dots represent the Hardy-Weinberg expectation.

356

## 357 References

- 358 1. Raven, P., Fauquet, C., Swaminathan, M.S., Borlaug, N. & Samper, C. Where Next  
359 for Genome Sequencing? *Science* **311**, 468-468 (2006).
- 360 2. Olsen, K.M. & Schaal, B.A. Evidence on the origin of cassava: Phylogeography of  
361 *Manihot esculenta*. *Proceedings of the National Academy of Sciences* **96**, 5586-  
362 5591 (1999).
- 363 3. Wang, W. et al. Cassava genome from a wild ancestor to cultivated varieties. *Nat*  
364 *Commun* **5**, (2014).
- 365 4. McDonald, M.J., Rice, D.P. & Desai, M.M. Sex speeds adaptation by altering the  
366 dynamics of molecular evolution. *Nature* **531**, 233-236 (2016).
- 367 5. McKey, D., Elias, M., Pujol, B. & Duputié, A. The evolutionary ecology of clonally  
368 propagated domesticated plants. *New Phytologist* **186**, 318-332 (2010).
- 369 6. Charlesworth, D. & Willis, J.H. The genetics of inbreeding depression. *Nat Rev*  
370 *Genet* **10**, 783-796 (2009).
- 371 7. Rojas, M.C. et al. Analysis of Inbreeding Depression in Eight S1 Cassava Families.  
372 *Crop Science* **49**, 543-548 (2009).
- 373 8. Nuwamanya, E., Herselman, L. & Ferguson, M. Segregation of selected  
374 agronomic traits in six S1 cassava families. *Journal of Plant Breeding and Crop*  
375 *Science* **3**, 154-160 (2011).
- 376 9. Allem, A.C. The closest wild relatives of cassava ( *Manihot esculenta* Crantz).  
377 *Euphytica* **107**, 123-133 (1999).
- 378 10. Bredeson, J.V. et al. Sequencing wild and cultivated cassava and related species  
379 reveals extensive interspecific hybridization and genetic diversity. *Nat Biotech*  
380 **34**, 562-570 (2016).
- 381 11. Prochnik, S. et al. The Cassava Genome: Current Progress, Future Directions.  
382 *Tropical Plant Biology* **5**, 88-94 (2012).
- 383 12. Browning, Brian L. & Browning, Sharon R. Genotype Imputation with Millions of  
384 Reference Samples. *The American Journal of Human Genetics* **98**, 116-126.
- 385 13. Chia, J.-M. et al. Maize HapMap2 identifies extant variation from a genome in  
386 flux. *Nat Genet* **44**, 803-807 (2012).
- 387 14. Tennessen, J.A. et al. Evolution and Functional Impact of Rare Coding Variation  
388 from Deep Sequencing of Human Exomes. *Science* **337**, 64-69 (2012).
- 389 15. Fu, W. et al. Analysis of 6,515 exomes reveals the recent origin of most human  
390 protein-coding variants. *Nature* **493**, 216-220 (2013).
- 391 16. Davydov, E.V. et al. Identifying a High Fraction of the Human Genome to be  
392 under Selective Constraint Using GERP++. *PLoS Comput Biol* **6**, e1001025 (2010).
- 393 17. Rodgers-Melnick, E. et al. Recombination in diverse maize is stable, predictable,  
394 and associated with genetic load. *Proceedings of the National Academy of*  
395 *Sciences* **112**, 3823-3828 (2015).
- 396 18. Mackay, T.F.C. et al. The *Drosophila melanogaster* Genetic Reference Panel.  
397 *Nature* **482**, 173-178 (2012).
- 398 19. Kumar, P., Henikoff, S. & Ng, P.C. Predicting the effects of coding non-  
399 synonymous variants on protein function using the SIFT algorithm. *Nat. Protocols*  
400 **4**, 1073-1081 (2009).



- 401 20. Mezmouk, S. & Ross-Ibarra, J. The Pattern and Distribution of Deleterious  
402 Mutations in Maize. *G3: Genes/Genomes/Genetics* **4**, 163-171 (2014).
- 403 21. Horvath, P. & Barrangou, R. CRISPR/Cas, the Immune System of Bacteria and  
404 Archaea. *Science* **327**, 167-170 (2010).
- 405 22. Keller, P.J. & Knop, M. Evolution of Mutational Robustness in the Yeast Genome:  
406 A Link to Essential Genes and Meiotic Recombination Hotspots. *PLoS Genet* **5**,  
407 e1000533 (2009).
- 408 23. Hussin, J.G. et al. Recombination affects accumulation of damaging and disease-  
409 associated mutations in human populations. *Nat Genet* **47**, 400-404 (2015).
- 410 24. Haddrill, P.R., Halligan, D.L., Tomaras, D. & Charlesworth, B. Reduced efficacy of  
411 selection in regions of the Drosophila genome that lack crossing over. *Genome*  
412 *Biology* **8**, 1-9 (2007).
- 413 25. Jørgensen, K. et al. Cassava Plants with a Depleted Cyanogenic Glucoside  
414 Content in Leaves and Tubers. Distribution of Cyanogenic Glucosides, Their Site  
415 of Synthesis and Transport, and Blockage of the Biosynthesis by RNA  
416 Interference Technology. *Plant Physiology* **139**, 363-374 (2005).
- 417 26. Conn, E.E. Cyanogenic Compounds. *Annual Review of Plant Physiology* **31**, 433-  
418 451 (1980).
- 419 27. Chen, H., Patterson, N. & Reich, D. Population differentiation as a test for  
420 selective sweeps. *Genome Research* **20**, 393-402 (2010).
- 421 28. Marsden, C.D. et al. Bottlenecks and selective sweeps during domestication have  
422 increased deleterious genetic variation in dogs. *Proceedings of the National*  
423 *Academy of Sciences* **113**, 152-157 (2016).
- 424 29. Agrawal, A.F. & Whitlock, M.C. Inferences About the Distribution of Dominance  
425 Drawn From Yeast Gene Knockout Data. *Genetics* **187**, 553-566 (2011).
- 426 30. Hohenlohe, P.A. et al. Population Genomics of Parallel Adaptation in Threespine  
427 Stickleback using Sequenced RAD Tags. *PLoS Genet* **6**, e1000862 (2010).
- 428 31. Danecek, P. et al. The variant call format and VCFtools. *Bioinformatics* **27**, 2156-  
429 2158 (2011).
- 430 32. Bradbury, P.J. et al. TASSEL: software for association mapping of complex traits  
431 in diverse samples. *Bioinformatics* **23**, 2633-2635 (2007).
- 432 33. Beissinger, T.M., Rosa, G.J., Kaeppler, S.M., Gianola, D. & de Leon, N. Defining  
433 window-boundaries for genomic analyses using smoothing spline techniques.  
434 *Genetics Selection Evolution* **47**, 1-9 (2015).
- 435 34. Du, Z., Zhou, X., Ling, Y., Zhang, Z. & Su, Z. agriGO: a GO analysis toolkit for the  
436 agricultural community. *Nucleic Acids Research* **38**, W64-W70 (2010).
- 437 35. Supek, F., Bošnjak, M., Škunca, N. & Šmuc, T. REVIGO Summarizes and Visualizes  
438 Long Lists of Gene Ontology Terms. *PLoS ONE* **6**, e21800 (2011).
- 439 36. Henn, B.M. et al. Distance from sub-Saharan Africa predicts mutational load in  
440 diverse human genomes. *Proceedings of the National Academy of Sciences* **113**,  
441 E440-E449 (2016).
- 442
- 443