**Genome-wide association mapping of correlated traits in cassava: dry matter and total carotenoid content**

Ismail Y. Rabbi[1]*, Lovina I. Udoh[1], Marnin Wolfe[2], Elizabeth Y. Parkes[1], Melaku A. Gedil[1], Alfred Dixon[1], Punna Ramu[3], Jean-Luc Jannink[2,4], Peter Kulakow[1].

1. International Institute of Tropical Agriculture (IITA), PMB 5320 Ibadan, Oyo, Nigeria
2. Department of Plant Breeding and Genetics, Cornell University, Ithaca, NY, USA
3. Institute of Genomic Diversity, Cornell University, Ithaca, NY, USA.
4. USDA-ARS, R.W. Holley Center for Agriculture and Health, Ithaca, NY, USA

* Corresponding author, email: I.Rabbi@cgiar.org

**ABSTRACT**

**Cassava (*Manihot esculenta* (L.) Crantz) is a starchy root crop cultivated in the tropics for fresh consumption and commercial processing. Dry matter content and micronutrient density, particularly of provitamin A – traits that are negatively correlated – are among the primary selection objectives in cassava breeding. This study aimed at identifying genetic markers associated with these traits and uncovering the potential underlying cause of their negative correlation – whether linkage and/or pleiotropy. A genome-wide association mapping using 672 clones genotyped at 72,279 SNP loci was carried out. Root yellowness was used indirectly to assess variation in carotenoid content. Two major loci for root yellowness was identified on chromosome 1 at positions 24.1 and 30.5 Mbp. A single locus for dry matter content that co-located with the 24.1 Mbp peak for carotenoid content was identified. Haplotypes at these loci explained a large proportion of the phenotypic variability. Evidence of mega-base-scale linkage disequilibrium around the major loci of the two traits and detection of the major dry matter locus in independent analysis for the white- and yellow-root subpopulations suggests that physical linkage rather that pleiotropy is more likely to be the cause of the negative correlation between the target traits. Moreover, candidate genes for carotenoid (*phytoene synthase*) and starch biosynthesis (*UDP-glucose pyrophosphorylase* and *sucrose synthase*) occurred in the vicinity of the identified locus at 24.1 Mbp. These findings elucidate on the genetic architecture of carotenoids and dry matter in cassava and provides an opportunity to accelerate genetic improvement of these traits.**

**CORE IDEAS**
- Cassava, a starchy root crop, is a major source of dietary calories in the tropics.
- Most varieties consumed are poor in micronutrients, including pro-vitamin A.
- These two traits are governed by few major loci on chromosome one.
- Genetic linkage, rather than pleiotropy, is the most likely cause of their negative correlation.

**INTRODUCTION**

Cassava (*Manihot esculenta* (L.) Crantz) is one of the most important food and feed crops in the tropics and Africa accounts for more than half of the total world-wide production of

47  270.3 million tonnes (http://faostat3.fao.org/, accessed 26.03.2016). Because of its
48  remarkable tolerance to drought (El-Sharkawy, 1993), its ability to grow in poor soils (Cock,
49  1982), and its perennial nature which allows it to be harvested as and when required, this
50  heterozygous and clonally propagated species plays a particularly important role in food
51  security for millions of small-holder farmers in developing countries. Moreover, cassava is
52  increasingly being cultivated for commercial processing to convert its storage roots into
53  dehydrated chips, flour and starch (Balagopalan, 2002). Dry matter content, of which a large
54  proportion is starch, is therefore a primary factor that defines adoption of new cassava
55  varieties by farmers and the market value of harvested roots (Okechukwu and Dixon, 2008).
56  As a result, breeding of improved varieties with high dry matter content is one of the
57  primary objectives of cassava genetic improvement programs in the world.
58
59  Another important target trait for cassava improvement in developing countries is
60  biofortification for micronutrients (Pfeiffer and McClafferty, 2007; Saltzman et al., 2013).
61  Most varieties grown and consumed throughout the world have white storage roots with
62  negligible amounts of micronutrients in general, and provitamin A in particular (Welsch et
63  al., 2010). Dietary diversification and breeding of farmer-preferred improved varieties with
64  higher nutritional density are complementary approaches used in addressing potential
65  micronutrient deficiency associated with consumption of cassava as the sole staple food
66  (Sayre et al., 2011). The crop's gene-pool exhibits considerable natural variation for storage
67  root carotenoids that can be tapped for breeding of biofortified varieties, with some
68  breeding populations reported to accumulate as much as 25.8 µg/g fresh root weight
69  (Ceballos et al., 2013; Sánchez et al., 2014).
70
71  Despite availability of natural genetic diversity in the global germplasm that is relevant to
72  breeding for increased dry matter and total carotenoid contents, improving these traits
73  through phenotype-based recurrent selections is a lengthy process, due to the breeding
74  complexities associated with the species including an annual cropping cycle of 12 to 24
75  months and low multiplication rate of planting materials. Understanding the genetic basis of
76  variation in these traits is essential for increasing their selection efficiency and the rate of
77  genetic gain. More importantly, several studies using diverse germplasm have reported that
78  dry matter and carotenoid content are negatively correlated, with $r$ values ranging from -0.1
79  to -0.5 (Marín Colorado et al., 2009; Akinwale et al., 2010; Esuma et al., 2012; Ceballos et al.,
80  2013; Njoku et al., 2015). Despite its significant implication in breeding, the genetic basis of
81  this correlation – whether it is due to genetic linkage or pleiotropy – is not understood.
82
83  Several mapping studies using either Bulk Segregant Analysis (BSA) or Quantitative Trait Loci
84  (QTL) mapping of S1 or F1 populations have been reported separately for dry matter and
85  carotenoid content (Balyejusa Kizito et al., 2007; Marín Colorado et al., 2009; Welsch et al.,
86  2010; Morillo C et al., 2013; Njoku et al., 2014). The mapping resolution from single-cross
87  experimental populations is expected to be limited due to the use of sparse genetic maps
88  and the limited number of recombination events observed (Hamblin et al., 2011). Moreover,
89  QTLs from such bi-parental populations may not provide insight into the tremendous
90  genetic and phenotypic variation of the larger gene pool (Zhao et al., 2011). The increased
91  availability of genomic resources for cassava, including the chromosome-scale reference
92  genome and integrated linkage map (Prochnik et al., 2012; International Cassava Genetic
93  Map Consortium (ICGMC), 2014) and high-density genotyping using next-generation

94  sequencing (Rabbi et al., 2014a; b) makes it possible to use genome-wide association
95  (GWAS) mapping to dissect the phenotypic diversity of cassava germplasm with respect to
96  dry matter and carotenoid content. GWAS, which takes advantage of natural linkage
97  disequilibrium (LD) generated by ancestral mutation, drift, and recombination events in
98  diverse germplasm, offers the possibility to overcome the shortcomings of traditional bi-
99  parental QTL mapping. These advantages mean GWAS is able to reveal a broader spectrum
100 of trait-linked allelic variation and thus may provide the most useful markers for marker-
101 assisted selection (MAS). Indeed, GWAS has already been applied in other crops such as
102 maize to study the genetic architecture of carotenoid accumulation (Harjes et al., 2008;
103 Owens et al., 2014; Suwarno et al., 2015). In cassava, Esuma et al., (2016) carried out a
104 GWAS study using a panel of partial inbreds (S1 and S2 generation) produced from eight
105 clones. Using this limited number of parents, they reported a single genomic region on
106 Chromosome 1 that underlies the variation in total carotenoid content. However, no joint
107 association analysis examining carotenoids and dry matter content has hitherto been
108 reported. Here, we present the results of a GWAS using a collection of more than 650
109 cassava clones representing diverse African germplasm genotyped at high-density using
110 genotyping-by-sequencing (Elshire et al., 2011). The population was phenotyped in two
111 locations for three consecutive field seasons. The results of this study will be used to
112 develop efficient strategies to breed for high dry matter and provitamin A content varieties.
113
114 **METHODS**
115
116 **Germplasm:** The present work was carried out using the Tropical Manihot Selection (TMS)
117 cultivars developed at the International Institute of Tropical Agriculture (IITA) in Nigeria. This
118 population, also known as the Genetic Gain collection, consist of more than 650 advanced
119 breeding lines and key landraces selected over four decades from 1970 (Okechukwu and
120 Dixon, 2008; Ly et al., 2013). The pedigree of the collection is mainly composed of crosses
121 between germplasm from West Africa and early introductions of CMD-tolerant lines arising
122 from interspecific hybridization between *Manihot glaziovii* and cultivated cassava at the
123 Amani station in Tanzania (Hahn et al., 1980). The collection also includes hybrid germplasm
124 from Latin America (Wolfe et al., 2016).
125
126 **Locations and experimental design**: The Genetic Gain population was planted using an
127 incomplete block design with two checks per block and single row of either 5 or 10 plants
128 spaced at $1m^2$. Data used for this study was collected in the 2012-2013, 2013-2014, and
129 2014-2015 field seasons in Ibadan (7.40° N, 3.90° E) and Ubiaja (6.66° N, 6.38° E). The trials
130 are usually planted in June, at the onset of the raining season in South West Nigeria and
131 harvested in June of the following year.
132
133 **Assessment of dry matter content and yellow color intensity of storage roots:** Dry matter
134 content was assessed using the oven-drying method. Eight fully developed roots were
135 randomly selected from each plot, peeled, chipped and thoroughly mixed. For each sample,
136 100g was weighed and oven-dried for 48 hours at $104^o$C till constant weight was achieved.
137 The samples were then re-weighed and the dry matter content was expressed as the
138 percentage of dry weight relative to fresh weight.
139
140 Because of the well-established linear relationship between intensity of yellow color and

3

141 carotenoid content in cassava storage roots (Pearson's coefficient, $r$, ranges from 0.81 to
142 0.89), we used root yellowness as an indirect measure of carotenoid content (Iglesias et al.,
143 1997; Chávez et al., 2005; Marín Colorado et al., 2009; Akinwale et al., 2010; Sánchez et al.,
144 2014). The relative difference among clones in the Genetic Gain population was assessed
145 using two complementary methods. The first was a visual gradation of yellow color using a
146 standard color-chart starting from one (white) to 7 (deep yellow). Due to the potential
147 subjectivity inherent in visual color scores, we complemented the color-chart method
148 through the use of a Minolta CR-410® chromameter. Approximately 100g of grated samples
149 from freshly peeled roots were placed in transparent Nasco Whirl-Pak® sampling bags and
150 four chromameter measurements taken in different sections of the bag. We chose the
151 commission internationale de l'éclairage (CIELAB) method that records color values in a
152 three-dimensional color space, where the L* coordinate corresponds to a lightness
153 coordinate, and the a* coordinate corresponds either to red (positive values) or to green
154 (negative values). Of importance to this study was the b* coordinate, whose positive values
155 represents yellow while the negative values represent blue. The illuminant used was D65
156 and calibration was done each day with a white ceramic.

157

158 **SNP genotyping:** DNA was extracted as described in Rabbi et al. (2014) and Genotyping-by-
159 sequencing was carried out as described by (Elshire et al., 2011). DNAs from the Genetic
160 Gain individuals were digested individually with *ApeK*I, a methylation sensitive restriction
161 enzyme that recognizes a five base-pair sequence (GCWGC, where W is either A or T). The
162 GBS sequencing libraries consisting of 95-plex DNA samples each were prepared by ligating
163 the digested DNA to unique sample identifier barcodes (nucleotide adapters) followed by
164 standard PCR. Sequencing was performed using Illumina HiSeq2500. The sequenced reads
165 from different genotypes were de-convoluted using their unique barcodes and aligned to
166 version 6.0 of the cassava reference genome (www.phytozome.org/cassava) with the
167 Bowtie 2 (Langmead and Salzberg, 2012). SNPs were discovered using the GBS pipeline
168 Version 2 implemented in TASSEL software (Glaubitz et al., 2014) and converted to dosage
169 format (0 = homozygous reference, 1 = heterozygous, 2 = homozygous non-reference
170 alleles). Missing data were filtered as described in (Wolfe et al., 2016) and imputed with the
171 glmnet algorithm in R (http://cran.r-project.org/web/packages/glmnet/index.html) (Wong
172 et al., 2014).

173

174 **Phenotypic data analysis:**
175 The phenotypic data across two locations and three years was collapsed to single best linear
176 unbiased predictor (BLUP) values for each clone by fitting the following mixed linear model
177 with the *lme4* package in R:
178 $$y_{lij} = \mu + c_l + \beta_i + c_l * \beta_i + \varepsilon_{lij}$$
179 Here, $y_{lij}$ represents raw phenotypic observations, $\mu$ is the grand mean, $c_l$ is a random
180 effects term for clone with $c_l \sim N(0, \sigma_l^2)$, $\beta_i$ is a fixed effect for the combination of location
181 and year harvested, $C_l * \beta_i$ is a random effect for genotype-by-environment variance, and
182 $\varepsilon_{lij}$ is the residual variance, assumed to be random and distributed $N(0, \sigma_e^2)$. Broad-sense
183 heritability for dry matter content and yellow color intensity was calculated according to (Ly
184 et al., 2013). Genetic correlation among traits was also calculated from BLUP values.

185

186 **Population structure and Genome-Wide Association Analyses:** Inherent population
187 structure and cryptic relatedness can lead to spurious associations in GWAS (Astle and

4

188     Balding, 2009). To control for these confounding factors, three standard GWAS models were
189     compared: a simple one-way ANOVA model with no correction (naïve model); a general
190     linear model (GLM) with the first five PCs of the SNP matrix as covariates (GLM + 5PCs); and
191     a mixed-linear model (MLM) using the five PCs and marker-estimated kinship matrix (Yu et
192     al., 2006). The models correcting for kinship and 5 PCs had the lowest inflation-factors as
193     determined from quartile-quartile (QQ) plots and therefore the lowest false-discovery rate
194     (**Supplementary Figure 1**). The association analyses were implemented in TASSEL (Bradbury
195     et al., 2007; Zhang et al., 2010). Association test P-values were considered significant when
196     more extreme than the Bonferroni threshold (with experiment-wise type I error rate of
197     0.05).
198

199     The patterns and extent of linkage disequilibrium (LD) in a population not only determines
200     the obtainable resolution in association mapping studies (Hamblin et al., 2011) but also has
201     strong implication in the interpretation of association peaks. Therefore the level of LD decay
202     and the local patterns of LD along each chromosome were determined by calculating intra-
203     chromosomal pairwise squared correlation ($r^2$) using PLINK (Purcell et al., 2007).
204

205     **RESULTS**
206     **SNP genotyping:** A total of 72,279 genome-wide SNP markers were called for the 672
207     genetic gain individuals after filtering for minor allele frequency threshold of 0.005. The
208     high-density coverage of SNPs resulted in an average of 4015 markers per chromosome,
209     ranging from 3101 on chromosome 16 to 5880 on chromosome 1.
210

211     **Phenotypic variability:** We investigated the phenotypic variation in dry matter content as
212     well as carotenoid-based intensity of yellow root color using a visual color chart and
213     chromameter. The dry matter content varied widely in the Genetic Gain population ranging
214     from 8.4% to 45% (average 28.6%, Table 1). About two-thirds of the evaluated clones have
215     white storage roots while the remaining showed a range of yellow color suggesting varying
216     levels of carotenoid content. On average, the visual score was 1.7 and ranged from 1 (white)
217     to 7 (yellow). The average chromameter measure of yellow color intensity (b* value) was
218     20.8 and ranged from 11.1 (white) to 40.8 (yellow). Dry matter was approximately normally
219     distributed while chromameter b* values showed a bimodal distribution (**Figure
220     1**) in which the first peak (b* values from 10 to 20) is associated with the white clones while
221     the second peak (b* values from 20 to 40) is associated with the variations among the
222     yellow clones.
223     Broad-sense heritability was high for root yellow color ($H^2$ = 0.87 and 0.82 for color chart
224     and chromameter b* values, respectively) but moderate for dry matter content ($H^2$ = 0.51).
225     These values are within the range of heritabilities reported previously for these traits
226     (Balyejusa Kizito et al., 2007; Ceballos et al., 2013). The relative importance of genotype-by-
227     environment variance ($V_{GxE}$) compared to genotype ($V_G$) variance was measured by the ratio
228     $V_{GxE}/V_G$. For all traits, the genetic variance component was larger than the genotype-by-
229     environment interaction variance. The interaction is minimal for the yellow color
230     measurements (0.054 and 0.173 for color chart and chromameter b* values, respectively).
231     For dry matter content, we observed a slightly higher interaction ratio of 0.214.
232

233     The BLUPs for dry matter content and gradation of yellow color were negatively correlated
234     in our germplasm collection (Pearson's correlation coefficient, *r* = -0.59; P-value < 0.0001),

235  indicating that clones with higher carotenoid content are more likely to have low dry matter
236  content (**Figure 2**) which confirms previous findings in cassava (Marín Colorado et al., 2009;
237  Akinwale et al., 2010; Esuma et al., 2012; Njoku et al., 2015). On the other hand, we found a
238  positive association between dry matter content and color lightness (chromameter L* value,
239  *r* = 0.60, P-value < 0.0001). The two measures of yellow color (i.e. color chart and
240  chromameter b* axis)  were strongly correlated (*r* = 0.96, P-value < 0.0001).
241
242  **Population structure:** Analysis of population structure in 672 accessions genotyped across
243  72,279 SNPs using PCA detected subtle genetic differentiation in the genetic gain collection,
244  with the first 10 PCs explaining about 23 % of the genetic variation. The first two principal
245  components, which accounted for 8% of the genetic variation, revealed genetic
246  differentiation between white and yellow-root clones (**Figure 3**).
247
248  **Linkage disequilibrium:** Several regions of extensive mega-base-scale LD were discovered in
249  chromosomes 1, 4 and 10 as well as smaller regions in other chromosomes (**Figure 4**).
250  Excluding results from chromosomes with large LD blocks (i.e. chromosomes 1, 4 and 10),
251  we found that on average, LD drops almost to background levels ($r^2 < 0.1$) at around 2 Mb in
252  the Genetic Gain population (**Figure 5**).
253
254  **Population-wide GWAS:**
255  **Variation in carotenoid content estimated by root yellowness:** The MLM-based GWAS
256  analysis for yellow color in the storage root parenchyma using both the color chart and the
257  chromameter-based methods uncovered the same major association regions occurring at
258  24.1 and 30.5 Mbp of chromosome 1 (**Figure 6**). This is not surprising given the high
259  correlation between the two color assessment methods. The first major peak is tagged by
260  marker S1_24121306 for visual gradation of color (-$\log_{10}$(p-value) of 21.8) and marker
261  S1_24159585 (-$\log_{10}$(p-value) of 18.8) for chromameter b* value (Table 2). The second peak
262  was tagged by the same marker, S1_30543382, for both measures of color intensity and
263  occurred 6.5 Mbp away from the first peak (-$\log_{10}$(p-value) of 10.54 and 10.78, for color
264  chart and chromameter, respectively). All other SNPs between the two major regions were
265  not significant at the Bonferroni significance threshold (P = 6.92e-07) (**Figure 7**). The LD
266  between SNPs S1_24121306 and S1_30543382 was 0.3, suggesting moderate non-random
267  segregation of alleles of the two markers in these regions.
268
269  **Dry matter content:** Genetic variation in dry matter content was found to be associated
270  with a major locus occurring at 24.1 Mbp region of chromosome 1 and tagged by marker
271  S1_24121306 (–log10(p-value) of 11.73). Importantly, this locus for dry matter content co-
272  locates with one of the two peaks found to be associated with carotenoid content (**Figure
273  7**). The genomic co-location of the major loci for dry matter content and root yellowness
274  suggests either a strong physical linkage between the genes underlying these important
275  traits or a pleiotropic effect. Distinguishing between these two possible causes is important
276  in cassava improvement efforts that target both traits. We therefore attempted to unravel
277  the genetic cause of the observed association by: (i) exploring the underlying linkage
278  disequilibrium patterns in the QTL regions on chromosome 1; (ii) carrying out independent
279  association analysis for dry matter content within the white root and yellow root
280  subpopulations; and (iii) searching for plausible biological explanation by identifying
281  candidate genes for both traits in the target region.

282
283 Exploration of the LD landscape along Chromosome 1 uncovered a mega-base-scale region
284 of low recombination extending from 22Mb to 33 Mb surrounding the association peaks for
285 dry matter content and yellow color (**Figure 7**). This region was recently shown to coincide
286 with a large *Manihot glaziovii* introgression segment (Bredeson et al., 2016) that traces back
287 to early breeding for resistance to cassava mosaic and cassava brown streak viruses in the
288 1930's (Hahn et al., 1980). Clustering of the Genetic Gain population based on identity-by-
289 descent relationship (i.e. a measure of how many alleles at any marker in each of the two
290 samples came from the same ancestral chromosomes) calculated using only markers from
291 this extensive LD region (2150 SNPs from markers S1_21567540 to S1_34950326) revealed
292 at least two major groups of accessions (**Supplementary Figure 2**), indicating presence of
293 few major haplotypes associated with the LD blocks.
294
295 **GWAS for dry matter content in white root and yellow root subpopulation:** If the
296 phenotypic association between dry matter and carotenoid contents and the colocation of
297 their association signals (~ 24.1 Mbp region) is largely caused by physical linkage rather than
298 by pleiotropy, the major dry matter locus should be detectable in both white root and
299 yellow root germplasm when analyzed independently. We therefore split the Genetic Gain
300 dataset into white root (n=210) and yellow root (n=427) subpopulations and repeated the
301 GWAS analysis. Clones that were at the borderline between yellow-root and white-root
302 were excluded from these analyses. To mitigate the loss of power as a result of double-
303 fitting markers in the MLM model both as a fixed effect tested for association and as a
304 random effect as part of the kinship (Lippert et al., 2011; Listgarten et al., 2012), the MLM
305 analysis was carried out using a kinship matrix calculated excluding markers from
306 chromosome 1.
307
308 We recovered the major dry matter content association signal in both the white root and
309 yellow-root subpopulation (**Figure 8**). Though coinciding with the locus identified in the
310 population-wide GWAS, the association signal in the white subpopulation was much
311 broader, extending from 24 to 33 Mbp and generally overlaps with the broad LD region of
312 the chromosome 1 (**Figure 8**). On the contrary, association signal for the yellow
313 subpopulation was relatively narrow. Survey of the underlying LD pattern in the same
314 chromosome region for the yellow subpopulation showed a recombination spot.
315
316 **Selection sweep associated with breeding for yellow-root varieties:** To determine whether
317 the breeding for carotenoid content trait in the Genetic Gain germplasm resulted in a
318 selection sweep around the major QTL region, we quantified genome-wide nucleotide
319 variation in the yellow root subpopulation (n = 210) and the non-yellow subpopulation (n =
320 427). A sliding-window scan of expected heterozygosity ($\pi$) and Tajima's D detected a ~ 6
321 Mb region with decrease in nucleotide diversity in the yellow compared to white-root
322 subpopulation around the first major carotenoid locus site (~ 24.1 Mbp) relative to its
323 chromosomal neighborhood (**Figure 9**).
324
325 **Proportion of variance explained by markers QTL haplotypes.**
326 To determine predictive ability of the discovered loci for yellow color intensity and dry
327 matter content, we carried out a multiple linear regression analysis using the *lm* function in
328 R and considered the top markers for these traits as independent variables and the traits

329 measurements as the response variables. A model considering the two major peaks
330 associated with gradation of yellow color as assessed using a color chart (S1_24121306 and
331 S1_30543382) returned an adjusted squared correlation ($R^2$) of 0.81. For the measure of
332 continuous variation in intensity of yellow color using chromameter (b* value), the adjusted
333 $R^2$ from same genomic regions (S1_24159585 and S1_30543382) was 0.70 while that for dry
334 matter content was moderate ($R^2 = 0.37$). This finding suggests that the major loci on
335 Chromosome 1 would be useful in Marker Assisted Selection breeding in cassava. Single or
336 joint allelic substitution effects at the associated loci with respect to chromameter b* value,
337 color chart and dry matter content is shown in **Figure 10**.
338
339 **Candidate genes:** The first of the two genomic regions associated color intensity (tagged by
340 SNP S1_24159585) was found ~ 4.5 Kbp away from phytoene synthase 2 (*PSY*2) in the
341 cassava version 6 reference genome. The *PSY*2 enzyme, named Manes.01G124200 and
342 located at 24,155,070 bp, is involved in the first dedicated step of the carotenoid
343 biosynthesis pathway in cassava roots which converts geranylgeranyl diphosphate to
344 phytoene (Welsch et al., 2010). Presence of the null versus the functional *PSY*2 allele is
345 responsible for the qualitative color difference between the white and the yellow roots,
346 respectively (Welsch et al., 2010; Rabbi et al., 2014a). Our study suggests that allelic
347 variation associated with increases in enzyme activity could contribute to deeper yellow by
348 increasing the flux into the pathway. No known candidate genes were found in the vicinity
349 of the second significant association signal on chromosome 1 occurring at 30.5 Mbp.
350
351 For dry matter content, we found two particular genes that are pivotal in central carbon
352 metabolism in the vicinity of top SNP linked to that trait. The first is UDP-glucose
353 pyrophosphorylase (named Manes.01G123000 in the cassava reference genome). This gene
354 which occurs at 24.06 Mbp region, plays a key role in carbohydrate metabolism, and is
355 strongly associated with the yield production both in grains and root crops (Smith, 2008;
356 Zeeman et al., 2010). UDP-glucose pyrophosphorylase was recently found to be up-
357 regulated during bulking of cassava storage roots (Yang et al., 2011; Wang et al., 2016). The
358 second key carbohydrate metabolism gene was sucrose synthase (named
359 Manes.01G123800), which occurred in 24.14 Mbp region. Finding of these potential
360 candidate genes for carotenoid and carbohydrate biosynthesis strongly favors the possibility
361 that the association between these two traits is caused by physical linkage rather than
362 pleiotropy. This hypothesis warrants further investigation.
363
364 **DISCUSSION AND CONCLUSION**
365 The present study revealed that the genetic architecture for dry matter content and
366 intensity of yellow color resulting from carotenoid accumulation in cassava roots is
367 governed by few major loci on chromosome 1 and explains the large repeatability estimates,
368 particularly for yellow color. These findings expand on those from previous genetic mapping
369 efforts for dry matter and carotenoid content. Using a candidate gene mapping approach,
370 Welsch et al. (2010) reported that a SNP mutation in the *PSY2* gene, leading to amino-acid
371 substitution, differentiates white and yellow cassava storage roots. Similarly, a bi-parental
372 QTL mapping study that used two clones from the Genetic Gain collection (TMS-I961089A
373 and TMEB117) also uncovered a single QTL peak whose confidence interval encompassed
374 the same PSY2 gene (Rabbi et al., 2014a). The F1 progenies from the TMS-I961089A x
375 TMEB117 population, also genotyped using the GBS method, segregated at an

376  approximately 1:1 ratio for white versus light-yellow roots, suggesting that the yellow-root
377  parent was heterozygous for the functional allele at the PSY2 locus. Kizito et al. (2007)
378  reported a QTL for dry matter content in a bi-parental population genotyped using SSR
379  markers that also corresponds to this region on chromosome 1.  More recently, Esuma et
380  al., (2016) reported a single genomic region on Chromosome 1 underlies the variation in
381  total carotenoid content in eight S1 and S2 partially inbred families. This peak, around 24.66
382  Mbp, is close to our first locus tagged by SNP S1_24121306. However that study did not look
383  at genetic architecture for dry matter content.
384
385  While the amount of total carotenoids in the Genetic Gain collection was not directly
386  estimated, previous studies of diverse cassava germplasm have consistently reported a
387  strong linear relationship between yellow color and carotenoid content (Pearson's
388  coefficient, $r$, ranging from 0.81 to 0.89) (Iglesias et al., 1997; Chávez et al., 2005; Marín
389  Colorado et al., 2009; Akinwale et al., 2010; Sánchez et al., 2014). Hence the results
390  obtained here should be useful for breeding efforts targeting breeding for improved
391  carotenoid content. Nevertheless, we propose to quantify total carotenoids and its
392  constituents as a future study to corroborate the current findings.
393
394  Given the importance of dry matter content in cassava, and the fact that we found a single
395  genomic region associated with this trait, further studies are warranted to fine-map and
396  validate the identity of the causal locus. To do this effectively would require different
397  populations that are lacking the wild introgression segments in chromosome 1. This will lead
398  to reduced LD and allow higher mapping resolution. Additionally, special crosses such as
399  nested-association mapping population design (Yu et al., 2008) using strategically selected
400  sets of parents will reduce the confounding effect of population structure. Given our marker
401  density and sample size, this study is sufficiently powered to find large effect alleles that are
402  common in the studied germplasm. To detect more QTLs of small effects will require a
403  larger association panel genotyped at higher density.
404
405  The use of a broad cassava diversity panel in GWAS not only provides the foundation to map
406  genomic regions associated with natural variation in dry matter and carotenoid content but
407  also allows us to unravel the genetic cause of the negative correlation between these traits,
408  that is, pleiotropy versus genetic linkage. In the context of breeding to simultaneously
409  increase carotenoid and dry matter content, the observed negative association between
410  these traits in our germplasm is undesirable. Several lines of investigation pointed to a
411  possibility of genetic linkage rather than pleiotropy to be the cause of the observed
412  association. Firstly, the genomic region harboring the QTLs for yellow color and dry matter
413  content was found to occur in chromosomal segments that exhibits low overall
414  recombination in this region compared to the genome-wide patterns. Recent work by
415  Bredeson et al. (2016) has shown that this chromosome 1 region harbors a large *M. glaziovii*
416  introgression that commonly occurs in the Genetic Gain collection. Secondly, independent
417  association analysis for dry matter content on the white and the yellow subpopulations
418  detected the same association signal although the QTL in the white subpopulation was
419  broader suggesting that the favorable alleles were located in non-recombining haplotype.
420  Thirdly, strong candidate genes for dry matter (UDP-glucose pyrophosphorylase and
421  sucrose synthase) and carotenoid content (phytoene synthase) were found in the vicinity
422  of the major association region (24.1 Mbp) of chromosome 1. Presence of these genes hints

423  at possibly distinct biological causes of the observed associations with the two traits. These
424  hypotheses need to be tested through functional genetics studies at these candidate genes.
425  Taken together, these findings suggest that the phenotypic correlation between dry matter
426  and carotenoid content is mainly caused by physical linkage of loci underlying these trails.
427  Moreover, Ortiz et al. (2011) found a fairly large positive correlation (r = 0.62) between
428  these traits. It is therefore possible that the nature of association (whether positive or
429  negative) is dependent on the allelic status at the linked dry matter and carotenoid
430  biosynthesis genes. We also detected a reduction of expected heterozygosity (π) around the
431  major gene region in the yellow versus white sub population. This suggests that the genetic
432  base for sources of favorable alleles with respect to carotenoid biosynthesis at this locus is
433  narrow, possibly arising from a single haplotype, which could be linked in *cis* to low-dry
434  matter alleles in the dry matter locus. Alternatively, balancing selection of the *M. glaziovii*
435  introgression in the white cassava sub population might be the cause of the higher levels of
436  heterozygosity relative to the yellow sub population.
437
438  Although cassava is a predominantly outcrossing species, its clonally propagated nature
439  means that modern varieties have undergone relatively few recombination cycles compared
440  to seed crops. Most accessions in the Genetic Gain collection are not far removed from
441  founder clones. Accordingly, the extent of LD in this study (~ 2 Mbp) is much greater than
442  the LD in maize (< 10 Kb) (Yan et al., 2009) as well as in grape (< 10 Kb) (Myles et al., 2011),
443  another clonal species. Moreover, the overall recombination pattern is far from
444  homogeneous owing to the persistent introgressions of *M. glaziovii* chromosomal segments
445  that are legacies of the historical breeding program in East Africa (Hahn et al., 1980;
446  Jennings, 1994). From these results, it is expected that the mapping resolution will vary
447  widely across the cassava genome depending mainly on whether a locus-of-interest occurs
448  in or outside the large-LD blocks.
449
450  This study presents a significant progress toward dissecting the genetic architecture of two
451  key breeding goal traits in cassava. The major loci associated with carotenoid content
452  variation and a single locus associated with dry matter content represents markers that will
453  be useful for marker-assisted selection in these traits. Although the results of the present
454  study suggests genetic linkage is more likely to be responsible for the negative correlation
455  between the studied traits, there is need for further investigations to confirm or reject this
456  hypothesis. For example, will dry matter content be increased by knocking out the PSY2
457  gene using gene-silencing methods (Lu et al., 2003; Burch-Smith et al., 2004; Fofana et al.,
458  2004)?. Alternately, could the activation of PSY2 in clones with high dry matter content and
459  lacking in carotenoids using gene-editing technologies like CRISPR-CAS9 (Hsu et al., 2014;
460  Sander and Joung, 2014) lead to not only carotenoid production and accumulation but also
461  lowering of dry matter content?
462

470    Kingdom Department for International Development.
471

## References

473    Akinwale, M.G.. b, R.D.. Aladesanwa, B.O.. Akinyele, A.G.O.. Dixon, and A.C.. Odiyi. 2010.
474        Inheritance of β-carotene in cassava (Manihot esculenta crantza). Int. J. Genet. Mol.
475        Biol. 2(10): 198–201.
476    Astle, W., and D. Balding. 2009. Population Structure and Cryptic Relatedness in Genetic
477        Association Studies. Stat. Sci. 24(4): 451–471.
478    Balagopalan, C. 2002. Cassava: biology, production and utilization (RJ Hillocks and JM
479        Thresh, Eds.). CABI, Wallingford.
480    Balyejusa Kizito, E., A.-C. Rönnberg-Wästljung, T. Egwang, U. Gullberg, M. Fregene, and A.
481        Westerbergh. 2007. Quantitative trait loci controlling cyanogenic glucoside and dry
482        matter content in cassava (Manihot esculenta Crantz) roots. Hereditas 144(4): 129–
483        136.
484    Bradbury, P.J., Z. Zhang, D.E. Kroon, T.M. Casstevens, Y. Ramdoss, and E.S. Buckler. 2007.
485        TASSEL: software for association mapping of complex traits in diverse samples.
486        Bioinformatics 23(19): 2633–2635.
487    Bredeson, J. V, J.B. Lyons, S.E. Prochnik, G.A. Wu, C.M. Ha, E. Edsinger-Gonzales, J.
488        Grimwood, J. Schmutz, I.Y. Rabbi, C. Egesi, P. Nauluvula, V. Lebot, J. Ndunguru, G.
489        Mkamilo, R.S. Bart, T.L. Setter, R.M. Gleadow, P. Kulakow, M.E. Ferguson, S. Rounsley,
490        and D.S. Rokhsar. 2016. Sequencing wild and cultivated cassava and related species
491        reveals extensive interspecific hybridization and genetic diversity. Nat. Biotechnol.
492        34(5): 562–570.
493    Burch-Smith, T.M., J.C. Anderson, G.B. Martin, and S.P. Dinesh-Kumar. 2004. Applications
494        and advantages of virus-induced gene silencing for gene function studies in plants.
495        Plant J. 39(5): 734–746.
496    Ceballos, H., N. Morante, T. Sánchez, D. Ortiz, I. Aragón, A.L. Chávez, M. Pizarro, F. Calle, and
497        D. Dufour. 2013. Rapid Cycling Recurrent Selection for Increased Carotenoids Content
498        in Cassava Roots. Crop Sci. 53(6): 2342.
499    Chávez, a. L., T. Sánchez, G. Jaramillo, J.M. Bedoya, J. Echeverry, E. a. Bolaños, H. Ceballos,
500        and C. a. Iglesias. 2005. Variation of quality traits in cassava roots evaluated in
501        landraces and improved clones. Euphytica 143(1–2): 125–133.
502    Cock, J.H. 1982. Cassava: a basic energy source in the tropics. Science 218(4574): 755–762.
503    El-Sharkawy, M.A. 1993. Drought-tolerant cassava for Africa, Asia, and Latin America.
504        Bioscience 43(7): 441–451.
505    Elshire, R.J., J.C. Glaubitz, Q. Sun, J. a Poland, K. Kawamoto, E.S. Buckler, and S.E. Mitchell.
506        2011. A robust, simple genotyping-by-sequencing (GBS) approach for high diversity
507        species. PLoS One 6(5): e19379.
508    Esuma, W., L. Herselman, M.T. Labuschagne, P. Ramu, F. Lu, Y. Baguma, E.S. Buckler, and
509        R.S. Kawuki. 2016. Genome-wide association mapping of provitamin A carotenoid
510        content in cassava. Euphytica.
511    Esuma, W., P. Rubaihayo, A. Pariyo, R. Kawuki, B. Wanjala, I. Nzuki, J.J. Harvey, and Y.
512        Baguma. 2012. Genetic Diversity of Provitamin A Cassava in Uganda. J. Plant Stud. 1(1):
513        60–71.
514    Fofana, I.B.F., A. Sangaré, R. Collier, C. Taylor, and C.M. Fauquet. 2004. A geminivirus-
515        induced gene silencing system for gene function validation in cassava. Plant Mol. Biol.
516        56(4): 613–24.

517 Glaubitz, J.C., T.M. Casstevens, F. Lu, J. Harriman, R.J. Elshire, Q. Sun, and E.S. Buckler. 2014.
518     TASSEL-GBS: A high capacity genotyping by sequencing analysis pipeline. PLoS One 9(2).
519 Hahn, S.K., E.R. Terry, and K. Leuschner. 1980. Breeding cassava for resistance to cassava
520     mosaic disease. Euphytica 29(3): 673–683.
521 Hamblin, M.T., E.S. Buckler, and J.-L. Jannink. 2011. Population genetics of genomics-based
522     crop improvement methods. Trends Genet. 27(3): 98–106.
523 Harjes, C.E., T.R. Rocheford, L. Bai, T.P. Brutnell, C.B. Kandianis, S.G. Sowinski, A.E. Stapleton,
524     R. Vallabhaneni, M. Williams, E.T. Wurtzel, J. Yan, and E.S. Buckler. 2008. Natural
525     genetic variation in lycopene epsilon cyclase tapped for maize biofortification. Science
526     319: 330–333.
527 Hsu, P.D., E.S. Lander, and F. Zhang. 2014. Development and applications of CRISPR-Cas9 for
528     genome engineering. Cell 157(6): 1262–1278.
529 Iglesias, C., J. Mayer, L. Chavez, and F. Calle. 1997. Genetic potential and stability of
530     carotene content in cassava roots. Euphytica 94(3): 367–373.
531 International Cassava Genetic Map Consortium (ICGMC). 2014. High-resolution linkage map
532     and chromosome-scale genome assembly for cassava (Manihot esculenta Crantz) from
533     10 populations. G3 (Bethesda). 5(1): 133–44.
534 Jennings, D.L. 1994. Breeding for resistance to African cassava mosaic geminivirus in East
535     Africa. Trop. Sci. 34(1): 110–122.
536 Langmead, B., and S.L. Salzberg. 2012. Fast gapped-read alignment with Bowtie 2. Nat
537     Methods 9(4): 357–359.
538 Lippert, C., J. Listgarten, Y. Liu, C.M. Kadie, R.I. Davidson, and D. Heckerman. 2011. FaST
539     linear mixed models for genome-wide association studies. Nat. Methods 8(10): 833–
540     837.
541 Listgarten, J., C. Lippert, C.M. Kadie, R.I. Davidson, E. Eskin, and D. Heckerman. 2012.
542     Improved linear mixed models for genome-wide association studies. Nat. Methods
543     9(6): 525–526.
544 Lu, R., A.M. Martin-Hernandez, J.R. Peart, I. Malcuit, and D.C. Baulcombe. 2003. Virus-
545     induced gene silencing in plants. Methods 30(4): 296–303.
546 Ly, D., M. Hamblin, I. Rabbi, G. Melaku, M. Bakare, H.G. Gauch, R. Okechukwu, A.G.O. Dixon,
547     P. Kulakow, and J.-L. Jannink. 2013. Relatedness and Genotype × Environment
548     Interaction Affect Prediction Accuracies in Genomic Selection: A Study in Cassava. Crop
549     Sci. 53(4): 1312–1325.
550 Marín Colorado, J.A., H. Ramírez, and M. Fregene. 2009. Genetic mapping and QTL analysis
551     for carotenes in a S1 population of cassava. Acta Agron. Univ. Nac. Colomb. 58(1): 15–
552     21.
553 Morillo C, A.C., Y. Morillo C, and H. Ceballos L. 2013. Identification of QTLs for carotene
554     content in the genome of cassava (Manihot esculenta Crantz) and S1 population
555     validation. Acta Agronómica, Univ. Nac. Colomb. 62(3): 196–206.
556 Myles, S., A.R. Boyko, C.L. Owens, P.J. Brown, F. Grassi, M.K. Aradhya, B. Prins, A. Reynolds,
557     J.-M. Chia, D. Ware, C.D. Bustamante, and E.S. Buckler. 2011. Genetic structure and
558     domestication history of the grape. Proc. Natl. Acad. Sci. U. S. A. 108(9): 3530–3535.
559 Njoku, D.N., V.E. Gracen, S.K. Offei, I.K. Asante, E.Y. Danquah, C.N. Egesi, and E. Okogbenin.
560     2014. Molecular marke r analysis of F1 progenies and their parents for carotenoids i
561     nheritance in African cassava (Manihot esculenta Crantz). African J. Biotechnol. 13(40):
562     3999–4007.
563 Njoku, D.N., V.E. Gracen, S.K. Offei, I.K. Asante, C.N. Egesi, P. Kulakow, and H. Ceballos.

564      2015. Parent-offspring regression analysis for total carotenoids and some agronomic
565      traits in cassava. Euphytica 206(3): 657–666.

566 Okechukwu, R.U., and a. G.O. Dixon. 2008. Genetic Gains from 30 Years of Cassava Breeding
567      in Nigeria for Storage Root Yield and Disease Resistance in Elite Cassava Genotypes. J.
568      Crop Improv. 22(2): 181–208.

569 Ortiz, D., T. Sánchez, N. Morante, H. Ceballos, H. Pachón, M.C. Duque, A.L. Chávez, and A.F.
570      Escobar. 2011. Sampling strategies for proper quantification of carotenoid content in
571      cassava breeding. J. Plant Breed. Crop Sci. 3(1): 14–23.

572 Owens, B.F., A.E. Lipka, M. Magallanes-Lundback, T. Tiede, C.H. Diepenbrock, C.B. Kandianis,
573      E. Kim, J. Cepela, M. Mateos-Hernandez, C.R. Buell, E.S. Buckler, D. DellaPenna, M. a
574      Gore, and T. Rocheford. 2014. A Foundation for Provitamin A Biofortification of Maize:
575      Genome-Wide Association and Genomic Prediction Models of Carotenoid Levels.
576      Genetics 198(4): 1699–1716.

577 Pfeiffer, W.H., and B. McClafferty. 2007. HarvestPlus: Breeding Crops for Better Nutrition.
578      Crop Sci. 47(Supplement_3): S-88.

579 Prochnik, S., P.R. Marri, B. Desany, P.D. Rabinowicz, C. Kodira, M. Mohiuddin, F. Rodriguez,
580      C. Fauquet, J. Tohme, T. Harkins, D.S. Rokhsar, and S. Rounsley. 2012. The Cassava
581      Genome: Current Progress, Future Directions. Trop. Plant Biol. 5(1): 88–94.

582 Purcell, S., B. Neale, K. Todd-Brown, L. Thomas, M. a R. Ferreira, D. Bender, J. Maller, P.
583      Sklar, P.I.W. de Bakker, M.J. Daly, and P.C. Sham. 2007. PLINK: a tool set for whole-
584      genome association and population-based linkage analyses. Am. J. Hum. Genet.
585      81(September): 559–575.

586 Rabbi, I., M. Hamblin, M. Gedil, P. Kulakow, M. Ferguson, A.S. Ikpan, D. Ly, and J.L. Jannink.
587      2014a. Genetic mapping using genotyping-by-sequencing in the clonally propagated
588      cassava. Crop Sci. 54(4): 1384–1396.

589 Rabbi, I.Y., M.T. Hamblin, P.L. Kumar, M. a Gedil, A.S. Ikpan, J.-L. Jannink, and P. a Kulakow.
590      2014b. High-resolution mapping of resistance to cassava mosaic geminiviruses in
591      cassava using genotyping-by-sequencing and its implications for breeding. Virus Res.
592      186: 87–96.

593 Saltzman, A., E. Birol, H.E. Bouis, E. Boy, F.F. De Moura, Y. Islam, and W.H. Pfeiffer. 2013.
594      Biofortification : Progress toward a more nourishing future. Glob. Food Sec. 2(1): 9–17.

595 Sánchez, T., H. Ceballos, D. Dufour, D. Ortiz, N. Morante, F. Calle, T. Zum Felde, M.
596      Domínguez, and F. Davrieux. 2014. Prediction of carotenoids, cyanide and dry matter
597      contents in fresh cassava root using NIRS and Hunter color techniques. Food Chem.
598      151: 444–451.

599 Sander, J.D., and J.K. Joung. 2014. CRISPR-Cas systems for editing, regulating and targeting
600      genomes. Nat. Biotechnol. 32(4): 347–55.

601 Sayre, R., J.R. Beeching, E.B. Cahoon, C. Egesi, C. Fauquet, J. Fellman, M. Fregene, W.
602      Gruissem, S. Mallowa, M. Manary, B. Maziya-Dixon, A. Mbanaso, D.P. Schachtman, D.
603      Siritunga, N. Taylor, H. Vanderschuren, and P. Zhang. 2011. The BioCassava plus
604      program: biofortification of cassava for sub-Saharan Africa. Annu. Rev. Plant Biol. 62:
605      251–72.

606 Smith, A.M. 2008. Prospects for increasing starch and sucrose yields for bioethanol
607      production. Plant J. 54(4): 546–558.

608 Suwarno, W.B., K. V Pixley, N. Palacios-Rojas, S.M. Kaeppler, and R. Babu. 2015. Genome-
609      wide association analysis reveals new targets for carotenoid biofortification in maize.
610      Theor. Appl. Genet. 128(5): 851–864.

611  Wang, X., L. Chang, Z. Tong, D. Wang, Q. Yin, D. Wang, X. Jin, Q. Yang, L. Wang, Y. Sun, Q.
612      Huang, A. Guo, and M. Peng. 2016. Proteomics Profiling Reveals Carbohydrate
613      Metabolic Enzymes and 14-3-3 Proteins Play Important Roles for Starch Accumulation
614      during Cassava Root Tuberization. Sci. Rep. 6(January): 19643.
615  Welsch, R., J. Arango, C. Bär, B. Salazar, S. Al-Babili, J. Beltrán, P. Chavarriaga, H. Ceballos, J.
616      Tohme, and P. Beyer. 2010. Provitamin A accumulation in cassava (Manihot esculenta)
617      roots driven by a single nucleotide polymorphism in a phytoene synthase gene. Plant
618      Cell 22(10): 3348–56.
619  Wolfe, M.D., I.Y. Rabbi, C. Egesi, M. Hamblin, R. Kawuki, P. Kulakow, R. Lozano, D.P. Del
620      Carpio, P. Ramu, and J.-L. Jannink. 2016. Genome-Wide Association and Prediction
621      Reveals Genetic Architecture of Cassava Mosaic Disease Resistance and Prospects for
622      Rapid Genetic Improvement. Plant Genome 9(2): 1–248.
623  Wong, W.W.L., J. Griesman, and Z.Z. Feng. 2014. Imputing genotypes using regularized
624      generalized linear regression models. Stat. Appl. Genet. Mol. Biol. 13(5).
625  Yan, J., T. Shah, M.L. Warburton, E.S. Buckler, M.D. McMullen, and J. Crouch. 2009. Genetic
626      characterization and linkage disequilibrium estimation of a global maize collection
627      using SNP markers. PLoS One 4(12): e8451.
628  Yang, J., D. An, and P. Zhang. 2011. Expression Profiling of Cassava Storage Roots Reveals an
629      Active Process of Glycolysis/Gluconeogenesis. J. Integr. Plant Biol. 53(3): 193–211.
630  Yu, J., J.B. Holland, M.D. Mcmullen, and E.S. Buckler. 2008. Genetic Design and Statistical
631      Power of Nested Association Mapping in Maize. 551(January): 539–551.
632  Yu, J., G. Pressoir, W.H. Briggs, I. Vroh Bi, M. Yamasaki, J.F. Doebley, M.D. McMullen, B.S.
633      Gaut, D.M. Nielsen, J.B. Holland, S. Kresovich, and E.S. Buckler. 2006. A unified mixed-
634      model method for association mapping that accounts for multiple levels of relatedness.
635      Nat. Genet. 38(2): 203–208.
636  Zeeman, S.C., J. Kossmann, and A.M. Smith. 2010. Starch: its metabolism, evolution, and
637      biotechnological modification in plants. Annu. Rev. Plant Biol. 61: 209–234.
638  Zhang, Z., E. Ersoz, C.-Q. Lai, R.J. Todhunter, H.K. Tiwari, M. a Gore, P.J. Bradbury, J. Yu, D.K.
639      Arnett, J.M. Ordovas, and E.S. Buckler. 2010. Mixed linear model approach adapted for
640      genome-wide association studies. Nat. Genet. 42(4): 355–60.
641  Zhao, K., C.-W. Tung, G.C. Eizenga, M.H. Wright, M.L. Ali, A.H. Price, G.J. Norton, M.R. Islam,
642      A. Reynolds, J. Mezey, A.M. McClung, C.D. Bustamante, and S.R. McCouch. 2011.
643      Genome-wide association mapping reveals a rich genetic architecture of complex traits
644      in Oryza sativa. Nat. Commun. 2: 467.
645
646

14

**Tables**

**Table 1.** Summary of phenotype variation, variance components (±se) and broad-sense heritability ($H^2$) for dry matter, color chart and Chromameter CIELAB readings.

| Trait | DM (%) | TC-CHART | L* | a* | b* |
|---|---|---|---|---|---|
| Minimum | 8.4 | 1.0 | 69.5 | -3.0 | 11.1 |
| Average | 28.6 | 1.7 | 84.4 | -0.3 | 20.8 |
| Maximum | 45.4 | 7.0 | 90.2 | 4.6 | 40.8 |
| N | 3232 | 4237 | 1360 | 1360 | 1360 |
| $V_G$ | 16.52 (4.06) | 1.11 (1.05) | 5.22 (2.29) | 0.86 (0.92) | 16.52 (4.06) |
| $V_{GxE}$ | 3.53 (1.88) | 0.06 (0.24) | 0.77 (0.88) | 0.16 (0.41) | 2.87 (1.69) |
| $V_E$ | 1.91 (1.38) | 0.01 (0.11) | 0.53 (0.73) | 0.08 (0.28) | 1.11 (1.05) |
| Residual | 10.27 (3.20) | 0.17 (0.41) | 3.26 (1.80) | 0.30 (0.55) | 2.71 (1.65) |
| $H^2$ | 0.51 | 0.82 | 0.53 | 0.61 | 0.87 |

**Table 2.** Summary of significant associations between selected traits and SNP markers from the MLM analysis. Only results in the major loci from Chromosome 1 are shown.

| Trait | SNP | Chr | Position (bp) | P-value | Candidate genes and mid-position (bp) |
|---|---|---|---|---|---|
| Color chart | S1_24121306 | 1 | 24,121,306 | 1.74E-22 | Phytoene synthase (Manes.01G124200; 24,155,070 bp) |
| Color chart | S1_30543382 | 1 | 30,543,382 | 2.91E-11 | NA |
| Chromameter b* | S1_24159585 | 1 | 24,159,585 | 1.79E-19 | Phytoene synthase (Manes.01G124200; 24,155,070 bp) |
| Chromameter b* | S1_30543382 | 1 | 30,543,382 | 1.66E-11 | NA |
| Dry matter | S1_24121306 | 1 | 24,121,306 | 1.86E-12 | UDP-glucose pyrophosphorylase (Manes.01G123000; 24,061,652 bp); sucrose synthase (Manes.01G123800; 24,142,314 bp) |

Chr = Chromosome (version 6 of cassava reference genome);

15

## Figures



**Figure 1.** Distribution of phenotype for TCHART, Dry matter content, and chromameter (b*).
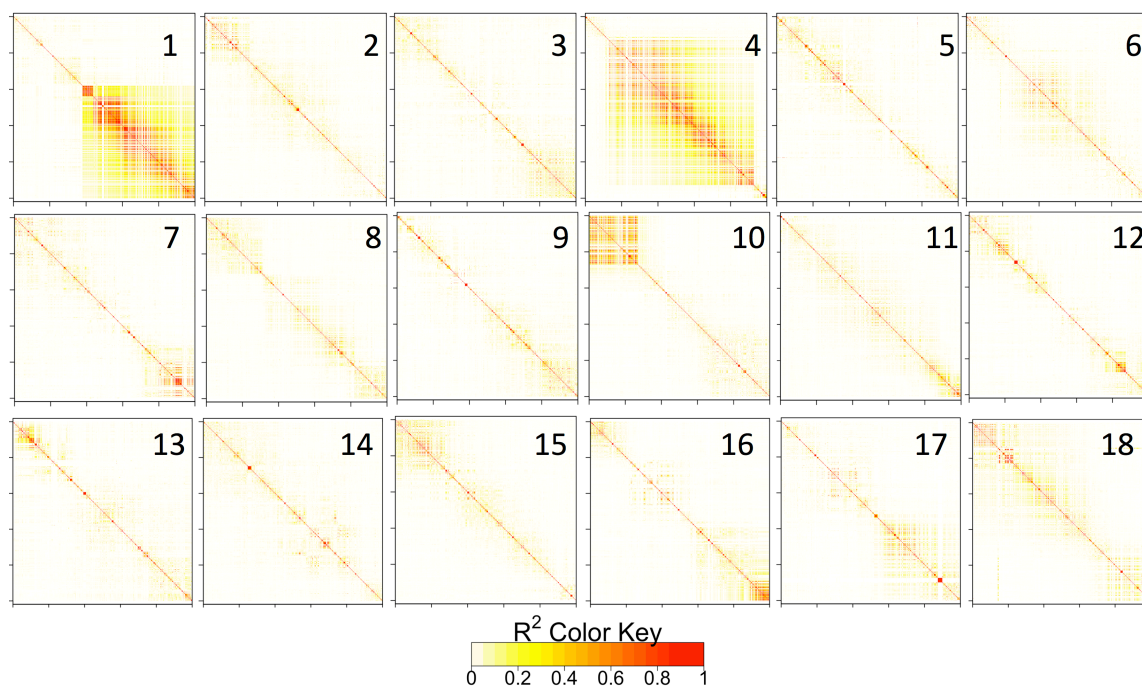


**Figure 2.** Relationship between dry matter content and root yellowness BLUPs (expressed as b* value of chromameter measurement). Different symbols denote the genotype at marker S1_24121306 that is associated with both dry matter content and root color intensity.
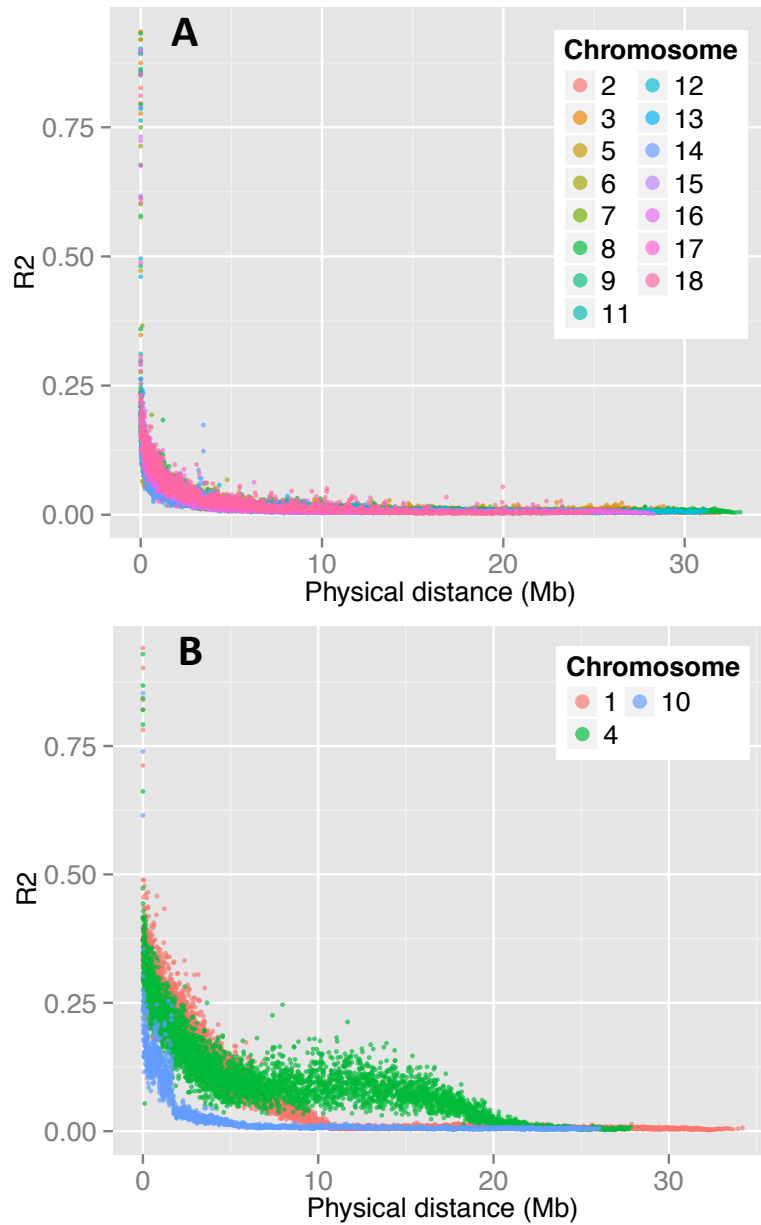
**Figure 3.** Population structure of the Genetic Gain collection. (A) PCA bi-plot of the first two axes; (B) Neighbor-joining dendrogram calculated from pairwise IBS distance. Yellow color highlights accessions with yellow roots.



**Figure 4.** Local pattern of linkage disequilibrium ($r^2$) along each of the 18 cassava chromosomes. Note the large LD blocks in chromosomes 1, 4 and 10. SNPs are arrayed according to their order, and not their physical position.

17

678
**Figure 5.** A Moving-average based LD decay profile in the Genetic Gain population. (A) all
chromosomes except 1, 4 and 10; (B) chromosomes 1, 4 and 10.

681

18

**Figure 6.** Genome-wide association results. Manhattan and Quantile-quantile plot of the MLM model for: root yellowness estimated using (A) chromameter b* value; and (B) color chart method; and (C) dry matter content. The red horizontal line indicates the genome-wide significance threshold.
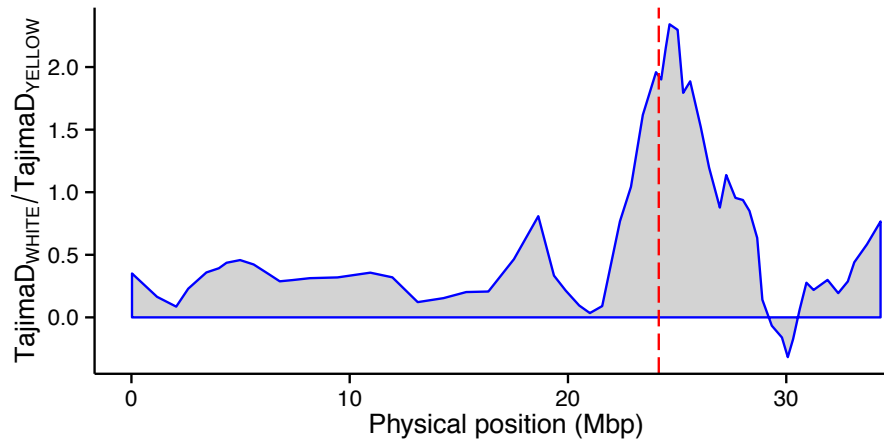
**Figure 7.** GWAS results for chromosome 1. Manhattan and Quantile-quantile plot of the MLM model for: (A) root yellowness as measured by chromameter b* value; (B) color chart; and (C) dry matter content. Note the common peak at ~ 24.1 Mbp region for the three traits. Red horizontal line indicates the genome-wide significance threshold. The SNPs are colored according to their degree of linkage disequilibrium ($r^2$) with the leading variant (i.e. top SNP for the first peak at 24.1 Mbp). The vertical blue lines in (A) and (B) denote the position of the carotenoid biosynthesis gene, phytoene synthase (24,155,070 bp), and those

696    on (C) denotes the positions of the UDP-glucose pyrophosphorylase (24,061,652 bp) and
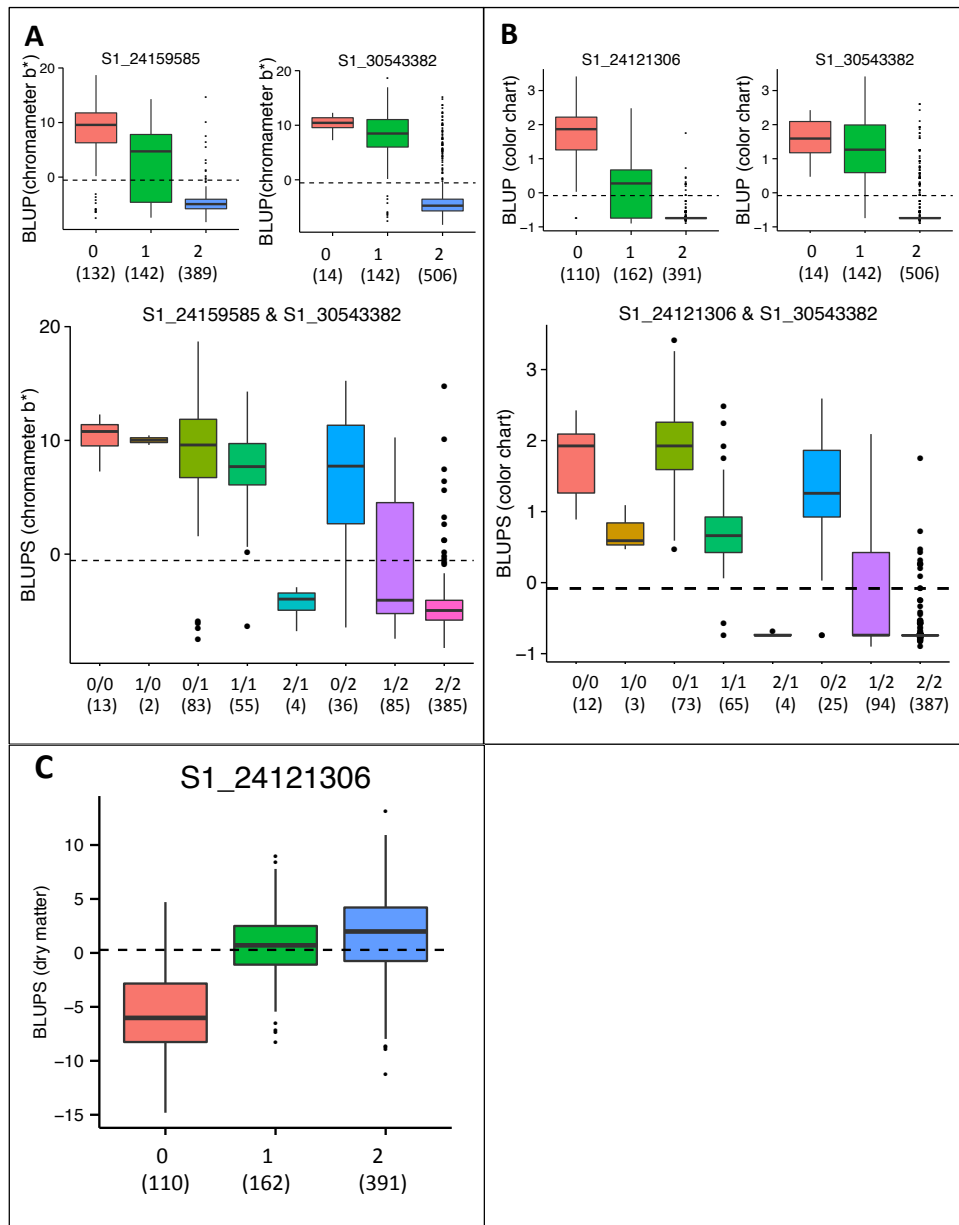697    sucrose synthase (24,142,314 bp).



698
699    **Figure 8.** Manhattan plots of the MLM analysis of the yellow root (top) and white root
700    subpopulations (bottom). Below each is an LD heatmap showing pairwise squared
701    correlation of alleles between markers along chromosome 1. Note the large number of SNPs
702    showing significant association with dry matter in the white subpopulation compared to
703    that of the yellow subpopulation. Red horizontal line indicates the genome-wide
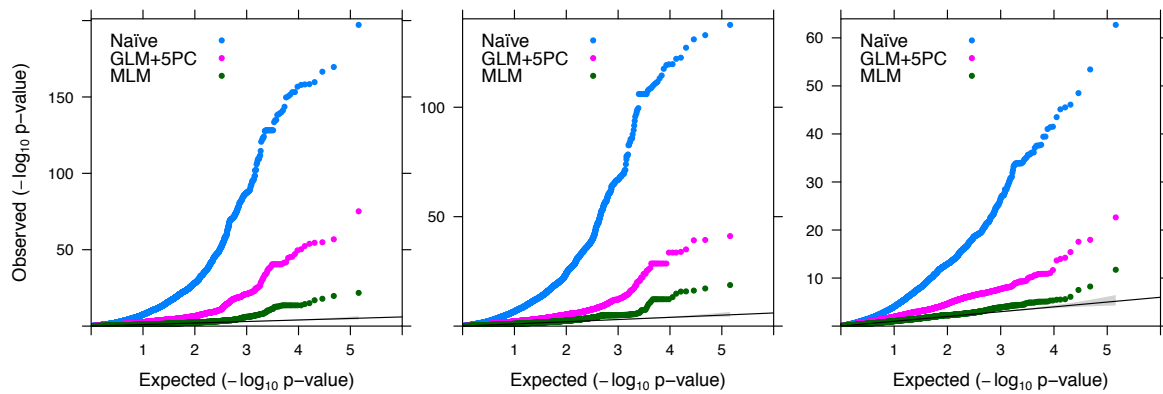704    significance threshold. The vertical blue lines are same as in Figure 7.
705

**Figure 9.** Selection sweep associated with positive selection for provitamin A trait varieties in the genetic gain population. Red dashed line indicates the position of SNP S1_24121306.
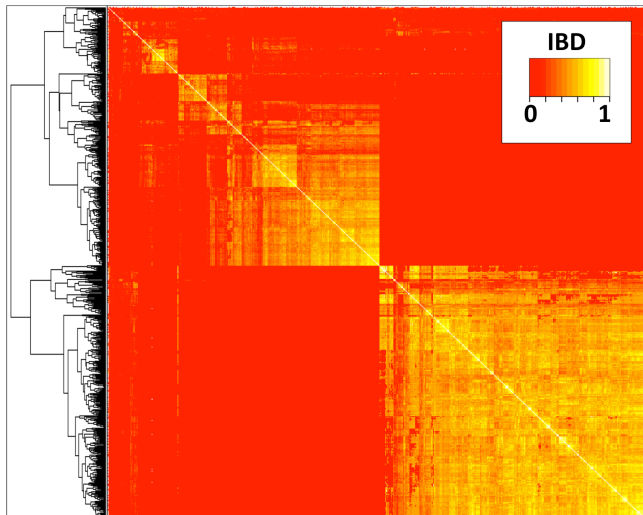
**Figure 10.** Effect of the most significantly associated markers on the BLUPs for yellow color measured by: (A) Chromameter (b* value) and (B) color chart in the Genetic Gain population. The boxplots show the effects of the most significantly associated SNPs at first and second peaks (above) and the two-locus haplotypes (below) on chromosomes 1. (C) Effect of the most significantly associated markers on the BLUPs for dry matter content. Alleles are coded as 0 = homozygous reference genome; 1 = heterozygous and 2 = homozygous non-reference genome. The dashed line represents the population mean of the BLUPs. The numbers in parenthesis below genotypic categories refer to the number of accessions for each genotype.



**Supplementary Figure 1.** Quantile–quantile plots for P-values obtained from simple GLM, GLM+5PCs and MLM model for color chart, chromameter b* value and dry matter content.



**Supplementary Figure 2.** Heatmap of identity-by-descent relationship using SNPs from large LD block in chromosome 1 around the major QTL region.

23