

1     Genome-wide Imputation Using the Practical Haplotype Graph in the Heterozygous  
2                                     Crop Cassava

3     Evan M Long\*, Peter J. Bradbury†‡, M. Cinta Romay†, Edward S. Buckler\*†‡, Kelly  
4     R Robbins\*

5             \* Plant Breeding and Genetics Section, School of Integrative Plant Science,  
6             Cornell University, Ithaca, NY 14853

7             † Institute for Genomic Diversity, Cornell University, Ithaca, NY 14853

8             ‡ United States Department of Agriculture-Agricultural Research Service, Robert  
9             W. Holley, Center for Agriculture and Health, Ithaca, NY 14853

10

11 Cassava Practical Haplotype Graph Imputation

12 Keywords: Cassava, Imputation, Haplotype, Practical Haplotype Graph, Genomic  
13 Prediction, Heterozygous, Beagle

14

15 Evan Long

16 175 Biotechnology Building

17 Ithaca, NY, 14853

18 503-413-0406

19

Eml255@cornell.edu

20

## ABSTRACT

21           Genomic applications such as genomic selection and genome-wide association  
22 have become increasingly common since the advent of genome sequencing. Genotype  
23 imputation makes it possible to infer whole genome information from limited input data, making  
24 large sampling for genomic applications more feasible, especially in non-model species where  
25 resources are less abundant. Imputation becomes increasingly difficult in heterozygous species  
26 where haplotypes must be phased. The Practical Haplotype Graph is a recently developed tool  
27 that can accurately impute genotypes, using a reference panel of haplotypes. The Practical  
28 Haplotype Graph is a haplotype database that implements a trellis graph to predict haplotypes  
29 using minimal input data. Genotyping information is aligned to the database and missing  
30 haplotypes are predicted from the most likely path through the graph. We showcase the ability  
31 of the Practical Haplotype Graph to impute genomic information in the highly heterozygous crop  
32 cassava (*Manihot esculenta*). Accurately phased haplotypes were sampled from runs of  
33 homozygosity across a diverse panel of individuals to populate the graph, which proved more  
34 accurate than relying on computational phasing methods. At 1X input sequence coverage, the  
35 Practical Haplotype Graph achieves a high concordance between predicted and true genotypes  
36 ( $R=0.84$ ), as compared to the standard imputation tool Beagle ( $R=0.69$ ). This improved  
37 accuracy was especially visible in the prediction of rare and heterozygous alleles. We validate  
38 the Practical Haplotype Graph as an accurate imputation tool in the heterozygous crop cassava,  
39 showing its potential for application in heterozygous species.

40

## INTRODUCTION

41           The past decade has seen an abundance of genomic sequence data produced  
42 for research and application in agricultural crops. With these new technologies, comes  
43 questions on how to effectively implement them (Torkamaneh *et al.* 2018). Two of the  
44 most common uses of genome-wide sequence data are genomic selection (GS) and

45 genome-wide association studies (GWAS). While most GWAS attempt to locate  
46 distinct, causative regions of the genome, GS incorporates all available markers to  
47 predict traits (Meuwissen *et al.* 2001). Genomic selection leverages a training set  
48 population that has both genotypic and phenotypic data to predict traits in a related  
49 germplasm with only genotypic data (Heffner *et al.* 2009). This allows breeders to both  
50 increase accuracy when selecting traits with low heritability and to accelerate the rate of  
51 selections by decreasing cycle time (Xu *et al.* 2020).

52         While sequencing data has become increasingly common in agricultural  
53 applications, the financial cost remains a challenge to widespread implementation.  
54 Reduced representation marker systems have been produced to limit costs of  
55 performing genomic analyses (Romay 2018), all of which vary in marker density and  
56 depth, cost, and genotype confidence. In scenarios with limited diversity, such as single  
57 breeding pools or post-bottleneck populations, individuals share large stretches of  
58 sequence. The strong association between alleles in these blocks, or their linkage  
59 disequilibrium (LD), determines the number and distribution of genotype markers  
60 needed to explain the genetic variation in the population. High density of markers  
61 becomes more important when performing analyses in populations where LD decays  
62 quickly as in species with high diversity or among unrelated individuals. High marker  
63 density can also be beneficial to incorporate knowledge on previously studied loci  
64 across the genome.

65         To affordably obtain high density genotypes or to bridge information between  
66 different marker platforms it becomes necessary to impute missing genotypes from  
67 available genotype data. Increasing the stability across genotyping platforms and

68 reducing per-sample costs becomes even more relevant in plant breeding scenarios,  
69 where many thousands of offspring are evaluated and changes in marker platform are  
70 common. Computational techniques to impute genome-wide information have been  
71 produced to bridge genotypic information from different marker panels and augment  
72 genotypic information from limited inputs (Yun *et al.* 2009). Genomic imputation  
73 methods often rely on a related training set with high confidence genotypic information  
74 to predict missing genotypes. These methods have been shown to improve consistency  
75 and efficiency of analyses of both genome wide associations (Spencer *et al.* 2009) and  
76 genomic selection (Cleveland *et al.* 2011).

77 Imputation is very common in genomic studies but is still plagued by barriers to  
78 high accuracy in many species. Known limitations of imputation stem from LD, allele  
79 frequencies, and population structure of the training population (Alipour *et al.* 2019).  
80 These difficulties are further compounded when working with a highly heterozygous  
81 crop, where both copies of the genome need to be modeled (Fragoso *et al.* 2016;  
82 Nazzicari *et al.* 2016). Heterozygosity introduces the challenge of phasing, or  
83 identifying which allele belongs to which copy of the genome, a challenge that is not  
84 limited to plants (Friedenberg and Meurs 2016). Imputation accuracy has been shown  
85 to affect the accuracy of genomic prediction in multiple scenarios (Pimentel *et al.* 2015;  
86 Wang *et al.* 2016; Van Den Berg *et al.* 2017). Highly accurate and less expensive  
87 imputation methods are needed to increase the gains made by GS by making  
88 genotyping more accurate and consistent. These improvements will enable research  
89 and breeding efforts to make accelerated gains, leading to more productive and  
90 adaptable crops in the changing global climate.

91           Rare variants contribute to the genetic load and overall performance of crops  
92 (Yang *et al.* 2017; Kremling *et al.* 2018; Kono *et al.* 2019), making high imputation  
93 accuracy, especially for alleles at low frequency, desirable for plant genomics  
94 applications. Diverse imputation tools exist and are often designed for different  
95 scenarios. One of the more common tools Beagle (Browning *et al.* 2018), which was  
96 designed for application in humans, works by leveraging LD between variants to predict  
97 missing genotypes. Beagle uses LD clustering to create an acyclic graph and a Hidden  
98 Markov model (HMM) to infer the most likely haplotype. Another method, EAGLE,  
99 leverages stretches of identity by descent (IBD) to perform long range phasing (Loh *et*  
100 *al.* 2016). In humans, where these imputation algorithms have been showcased, they  
101 have the advantage of large datasets with data from several thousands of individuals  
102 (Loh *et al.* 2016; Browning *et al.* 2018); this is not often possible in many plant breeding  
103 scenarios.

104           In maize, Beagle has difficulty accurately imputing rare variants, while a  
105 haplotype library based methods such as FILLIN can do so more easily (Swarts *et al.*  
106 2015). A recently developed method known as the Practical Haplotype Graph (PHG)  
107 was created to leverage known haplotypes in a graph structure to efficiently impute  
108 genotypes (Bradbury, In prep). The PHG simplifies the genome to a set of distinct  
109 regions of the genome, for which it defines haplotypes. These haplotypes are  
110 constructed from whole genome sequence data or genome assemblies and are used to  
111 construct a trellis graph, capturing the diversity of haplotypes at each range and the  
112 relationships between adjacent haplotype regions. Sequence reads are then aligned to  
113 the graph and a HMM is applied to predict the most likely haplotypes. By aligning reads

114 to pan-genome haplotypes, the PHG minimizes errors due to reference bias, poor  
115 alignment, and mis-called variants.

116 Here we showcase the potential application of the PHG in imputation of  
117 heterozygous crops. The PHG has already been shown to be an efficient tool for aiding  
118 imputation and genomic selection in breeding of the inbred cereal crop Sorghum  
119 (Jensen *et al.* 2020). It has also been implemented to impute genotypes in highly  
120 diverse maize lines (Franco *et al.* 2020). To show the utility of the PHG in a  
121 heterozygous crop we must overcome two distinct challenges: (1) obtaining phased  
122 haplotypes to populate the database and (2) modeling both copies of the genome  
123 accurately. Without an abundance of data, it is very difficult to obtain accurate phasing  
124 in a highly heterozygous species. This study will explore these challenges by imputing  
125 haplotypes from low-coverage skim sequencing, while comparing results to Beagle's  
126 performance.

127 To investigate the construction and performance of the PHG in a heterozygous  
128 scenario, we created a PHG for cassava (*Manihot esculenta*), a root crop with high  
129 levels of heterozygosity reinforced by centuries of clonal propagation. Cassava is a  
130 major caloric source for over half a billion people around the world, with a high  
131 concentration in sub-Saharan Africa (Parmar *et al.* 2017). Improved imputation in  
132 cassava could enable greater gains in breeding efforts to increase food security. In this  
133 study we utilize sequence data from the previously published HapMapII in cassava  
134 (Ramu *et al.* 2017), which includes whole genome sequence (WGS) data for 241  
135 cassava clones. This data is used to produce a PHG in cassava and illustrate its

136 effectiveness in genomic imputation in a heterozygous crop. We further validate these  
137 methods through genomic prediction and simulation.

138

## 139 MATERIALS AND METHODS

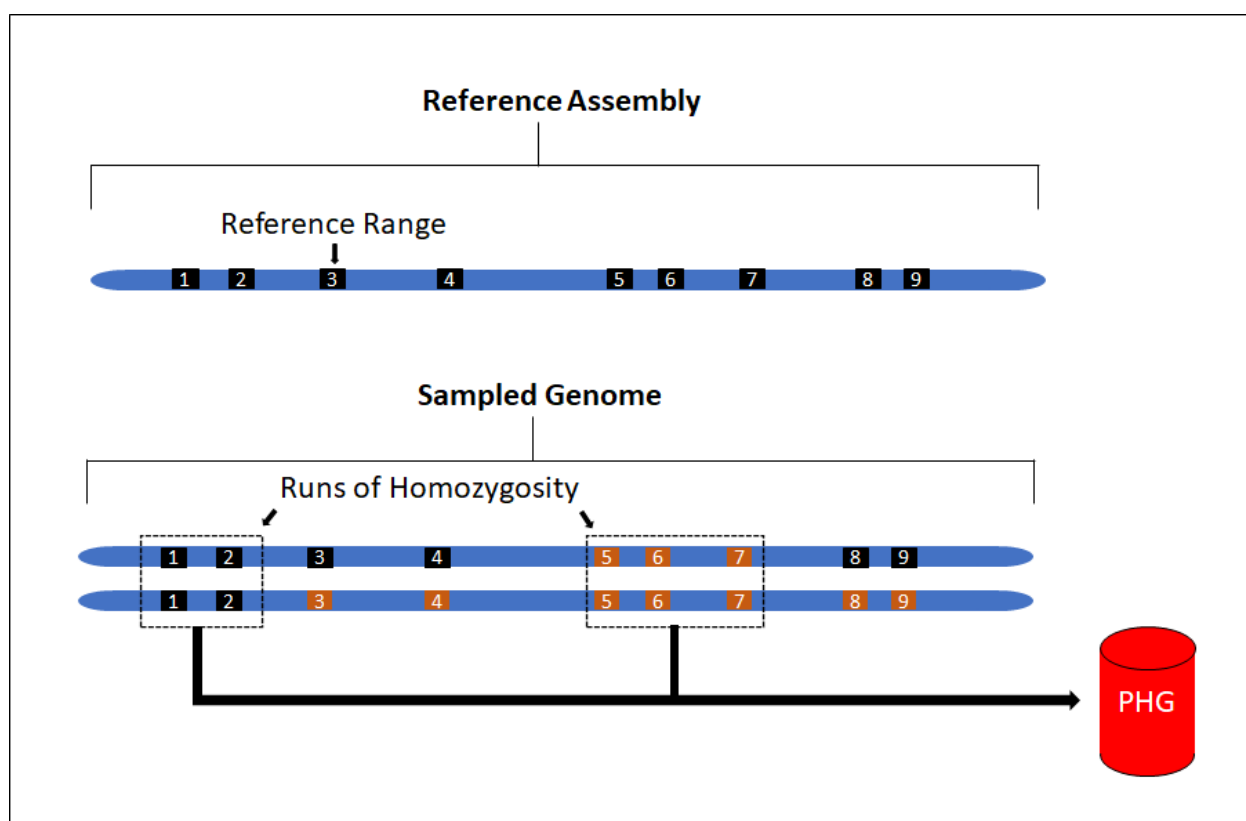
### 140 **Haplotype Sampling**

141 Genomic data was used from the second-generation Cassava Haplotype map  
142 consisting of 241 taxa, including both cultivated and wild germplasm (Ramu *et al.* 2017).  
143 Raw data is composed of short-read, whole genome sequence data from each taxon  
144 amounting to greater than 20X coverage on average. The high depth of the sequence  
145 data is necessary to accurately distinguish between heterozygous and homozygous  
146 variants. Haplotype regions, termed here as reference ranges, were defined by genic  
147 regions with additional 500bp flanking sequence from the Cassava V6 reference  
148 genome.

149 The detailed process of creating a PHG is outlined at  
150 [“https://bitbucket.org/bucklerlab/practicalhaplotypegraph/wiki/Home”](https://bitbucket.org/bucklerlab/practicalhaplotypegraph/wiki/Home) and has been  
151 described previously (Jensen *et al.* 2020; Franco *et al.* 2020). Here, we outline the  
152 specific steps taken to create a PHG in the heterozygous crop cassava. The major  
153 hurdle to producing a haplotype graph in a heterozygous species is obtaining accurately  
154 phased haplotypes. Because many of these cassava lines are cultivated taxa, we  
155 expected to find identical by descent (IBD) haplotypes brought about by generations of  
156 breeding within restricted breeding pools. These IBD segments provide confidently  
157 phased haplotypes as well as capturing their relationships to adjacent haplotypes (Fig.  
158 1). We identified and sampled these homozygous haplotypes which we inferred to



159 represent IBD haplotypes. This was done by measuring the number of heterozygous  
160 variants for each reference range in each taxon, then classifying those haplotypes as  
161 homozygous or not. The threshold for haplotypes to be considered IBD was determined  
162 empirically to be 0.001 heterozygous SNPs per base pair (Supplemental Fig. 1), as *de*  
163 *novo* mutations or errors in variant calling may produce low levels of perceived  
164 heterozygosity. This threshold was additionally validated by testing imputation accuracy  
165 of the PHG.



166  
167 **Figure 1. Haplotype view of the genome. Top) Representation of reference ranges**  
168 **informed from genic regions from the reference genome. Bottom) Haplotypes**  
169 **sampled from runs of homozygosity for use in PHG with different colors**  
170 **representing separate haplotypes at a given region (i.e. ranges 1,2,5,6,7 are**  
171 **homozygous and haplotypes can be sampled).**

172

173           After haplotypes were sampled from IBD regions of the genome, they were  
174 loaded as GVCF files into a PHG database. Similar haplotypes were then collapsed  
175 based on sequence similarity to produce a representative set of available haplotypes.  
176 Haplotypes are collapsed to make alignment more efficient, while retaining as much  
177 distinct haplotype information as possible. Collapsing is performed using an  
178 unweighted pair group method with arithmetic mean (upgma) tree from pairwise  
179 distance matrix from sequence variants to measure the similarity between haplotypes.  
180 Based on imputation accuracy tests, we chose a level of similarity (PHG parameter:  
181 maximum divergence) to collapse haplotypes of 0.001, corresponding to less than 1 in  
182 1000 nucleotide differences between haplotypes. This level of collapsing maintains  
183 high accuracy while collapsing redundant haplotypes (Supplemental Fig. 2). We then  
184 produced a pan-genome composed of consensus haplotypes representing the diversity  
185 of haplotypes.

## 186 **Predicting Haplotypes**

187           Once we obtained a set of consensus haplotypes, we implemented an HMM to  
188 infer genome-wide haplotypes from low depth genotyping data. Sparse genotype  
189 information was created by downsampling whole genome sequence data randomly  
190 using samtools to simulate skim sequencing. We randomly sampled 20 taxa from the  
191 cultivated varieties within the population to serve as a test set for downstream analyses.  
192 To test different levels of sequencing depth, we down-sampled reads to amounts  
193 estimated to represent 0.1X, 0.5X, 1X, 5X, and 10X single-end, whole genome  
194 sequence coverage.

195           These sampled sequences were aligned to the consensus haplotypes stored in  
196 the PHG to impute whole genome variants. A trellis graph is formed with every  
197 reference range representing separate ranges and the consensus haplotypes as nodes  
198 at each of those ranges. The most likely paths through the graph were then determined  
199 using an HMM Viterbi algorithm. Because cassava is heterozygous and diploid, this  
200 step produces the two most likely paths for each taxon. The emission and transition  
201 probability parameters of the HMM are defined by the genomes of the reference  
202 population used to build the database. The emission probabilities are calculated by  
203 considering the probability of two given haplotypes, given the aligned reads. The  
204 transition probabilities are defined by the edges between haplotypes in the PHG.

205           Due to the sparse sampling of IBD haplotypes from heterozygous taxa used to  
206 produce the PHG, the database lacked abundant transition information between  
207 adjacent reference ranges. To compensate for this, we aligned WGS for all 241 taxa  
208 used to create the database and predicted most likely paths through the graph. These  
209 paths were then used to augment the transition probabilities, without contributing any  
210 additional haplotypes.

### 211 **Beagle imputation**

212           We compared our imputation accuracy results to the common genotype  
213 imputation tool Beagle (Browning *et al.* 2018). Beagle was developed for the purpose of  
214 human data, but is a common tool used by many plant studies to impute missing  
215 genotypes. Because Beagle v4 can incorporate genotype likelihoods based on read  
216 depth, we used it for the imputation of the low depth sequence when it improved

217 accuracy, otherwise we utilized Beagle v5. We used the same HapMapII data from the  
218 241 clones to impute missing genotypes with Beagle.

## 219 **Genomic Prediction**

220 We used 57 clones from a single breeding program, to reduce effects of  
221 population structure, to determine the impact of imputation errors on genomic prediction  
222 accuracy using cross validation. Reads were down-sampled and imputed as previously  
223 described. Three root traits were used for genomic cross validation: fresh root yield,  
224 root size, and root number. Phenotypes for each clone were downloaded from  
225 CassavaBase.org, constituting 57 clones, spanning 23 years from 1996 to 2018, across  
226 13 locations in Africa. Ten-fold cross validation was performed by randomly selecting  
227 10% of the clones to hold out and predict using the remaining clones as a training set.  
228 The correlation between predicted phenotype and the observed BLUE was used as the  
229 prediction accuracy. We performed 50 replications as well as a single holdout  
230 prediction to measure genomic prediction accuracy. A single step model was  
231 performed:

$$232 \hat{y}_{ijkmp} = \mu + g_i + b_j + r_k + t_m + l_p + gt_{im} + gl_{ip}$$

233 Here,  $\hat{y}_{ijkmp}$  is the predicted trait and  $\mu$  is the fixed effect of the overall mean.  
234 Random effects were fitted as follows:  $g_i$  is genotype effect of the  $i^{\text{th}}$  clone,  $b_j$  is the  
235 effect of the  $j^{\text{th}}$  block,  $r_k$  is the effect of the  $k^{\text{th}}$  replicate,  $l_p$  is the location of the  $p^{\text{th}}$   
236 location,  $t_m$  is the effect of the  $m^{\text{th}}$  year,  $gl_{ip}$  is the interactive effect of the  $i^{\text{th}}$  clone and  
237 the  $p^{\text{th}}$  location, and  $gt_{im}$  is the interaction effect of the  $i^{\text{th}}$  clone and the  $m^{\text{th}}$  year. This  
238 was performed using the mixed model tool Echidna (Gilmour 2019).

239

240 The vectors of random effects for the model were distributed as

241  $\mathbf{g} \sim N(0, \mathbf{G}\sigma_g^2)$ ,  $\mathbf{b} \sim N(0, \mathbf{I}\sigma_b^2)$ ,  $\mathbf{r} \sim N(0, \mathbf{I}\sigma_r^2)$ ,  $\mathbf{t} \sim N(0, \mathbf{I}\sigma_t^2)$ ,  $\mathbf{l} \sim N(0, \mathbf{I}\sigma_l^2)$ ,  $\mathbf{gt} \sim N(0, \mathbf{I}\sigma_{gt}^2)$ ,  $\mathbf{gl} \sim N(0, \mathbf{I}\sigma_{gl}^2)$

242 Where  $\mathbf{G}$  is the genomic relationship matrix calculated using the “Eigenstrat” method of  
243 the R package SNPRelate (Zheng *et al.* 2012) and  $\mathbf{I}$  is the identity matrix.

### 244 **Pre-phased Haplotype PHG**

245 We investigated the viability of using computationally phased haplotypes to  
246 curate a PHG database rather than relying on IBD regions of the genome. First, we  
247 phased the variants from the 241 cassava clones using a combination of Beagle  
248 (Browning *et al.* 2018) and HAPCUT2 (Edge *et al.* 2017). These variants were used to  
249 create a PHG to be tested against the IBD version of the PHG. The second test utilized  
250 Oxford Nanopore (ONP) long-read sequencing from six cassava clones within the HMII  
251 population. High molecular weight DNA was extracted from young cassava leaves,  
252 selected for fragments 20-80 Kb long, and sequenced with MinION following the  
253 manufacturer recommendations. Variants were called using Guppy and their variants  
254 phased with WhatsHap (Schrinner *et al.* 2020). These six clones were then used to  
255 populate another PHG, we will identify as the “ONP6 PHG”. Larger reference ranges  
256 were divided into smaller regions to increase the probability of sampling correctly  
257 phased haplotypes. Twenty clones with the highest relationship to the six taxa with  
258 ONP data were used as the test set for these tests.

### 259 **Imputation from Simulated Genotypes**

260 A sample of 20 related individuals from the HapMapII population were selected  
261 to serve as parents for a simulated genotyping scenario. The genomes were phased  
262 using Beagle and then used to populate a PHG database. We then used these parents

263 to simulate 5 generations of random mating given a population size of 100  
264 (Supplemental Figure 3.). Forward genetic simulations were completed using SLiM  
265 (Haller and Messer 2019). Artificial short read-sequencing was then simulated for these  
266 offspring using neat-genreads (Stephens *et al.* 2016) at varied coverage levels. Reads  
267 were then aligned using bwa used to call and impute variants using Sentieon (Kendig *et*  
268 *al.* 2019) and Beagle. Reads were also aligned to the PHG formed from the original  
269 parents for imputation.

## 270 **Data Availability**

271 Supplementary files and scripts used for the production and testing of the cassava PHG  
272 can be found at [https://bitbucket.org/bucklerlab/p\\_cassava\\_phg](https://bitbucket.org/bucklerlab/p_cassava_phg). Genotype and  
273 phenotype data from HapMapII (Ramu *et al.* 2017) was downloaded from  
274 cassavabase.org. Support and methods for practical haplotype graph implementation  
275 can also be found at <https://bitbucket.org/bucklerlab/practicalhaplotypegraph/wiki/Home>.  
276 Raw Oxford nanopore sequence data for this project is available at NCBI BioProject ID  
277 PRJNA589272.

278

279

## RESULTS

### 280 **Haplotype Sampling**

281 To obtain phased haplotypes for the PHG we sampled haplotypes from  
282 homozygous regions of each clone. Centuries of clonal propagation and reported  
283 inbreeding depression (de Freitas *et al.* 2016) suggest cassava germplasm would be  
284 highly heterozygous, however, we found that, on average, ~20% of all reference ranges  
285 from each taxon were homozygous. This resulted in a high number of missing

286 haplotypes in each taxon, but a high confidence in the phased haplotypes that were  
287 sampled. Despite the variability in the number of homozygous samples by reference  
288 range, >90% of the reference ranges were homozygous in at least 10% of the HapMapII  
289 population (Supplemental Fig. 4). From these IBD haplotypes we were able to sample  
290 ~50% of the segregating sites. This proportion increased to 77% when considering  
291 sites with minor allele frequency above 5%, suggesting that many of the common  
292 variable sites have been sampled.

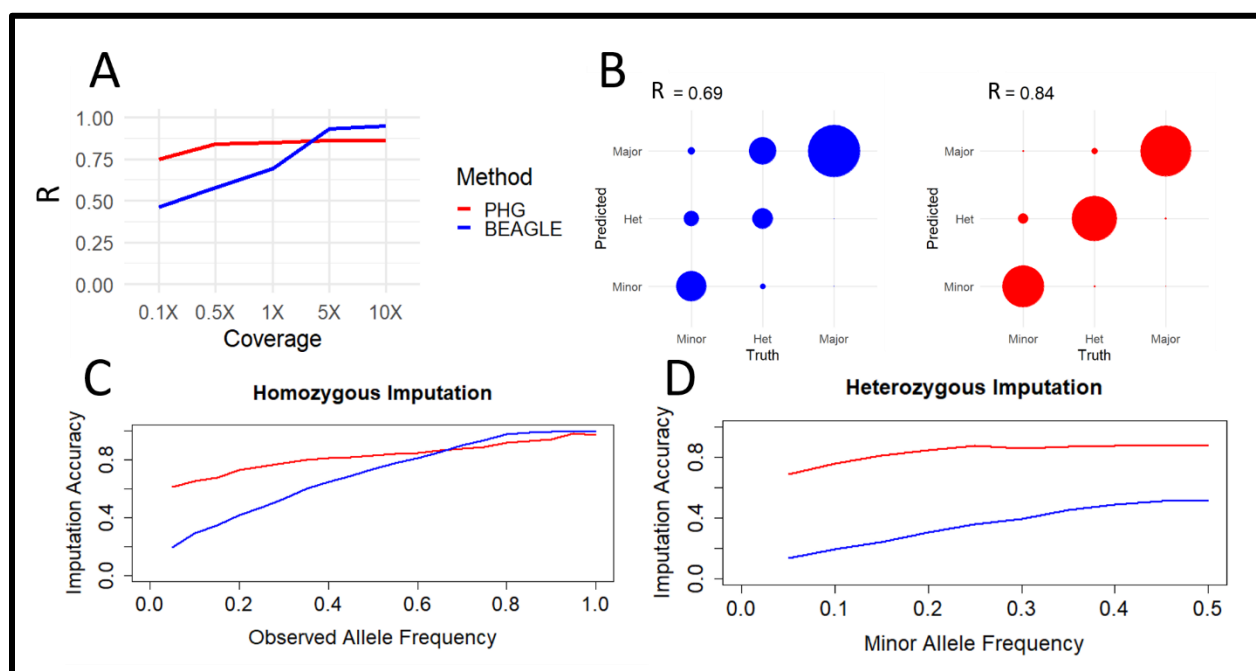
### 293 **Imputation and Genomic Prediction Accuracy**

294 Because imputation accuracy is dependent on the relative allele frequency and  
295 phase of the allele being called, we classified genotype calls by allele frequency class:  
296 homozygous major, homozygous minor, and heterozygous. In our analyses, imputation  
297 accuracy is defined as the ability of the imputation method to reconstitute genome-wide  
298 SNPs from the input data. We use the correlation between the predicted alleles and the  
299 true alleles (defined by HapMapII) as a metric to make the PHG and Beagle  
300 comparable, because the PHG utilizes reads and Beagle utilizes variants to make their  
301 predictions.

302 Imputation of skim sequence genotyping showed PHG methods had a large  
303 advantage over Beagle using low coverage sequence. At a level of 1X coverage  
304 random sequencing, the PHG predicted allele calls with a correlation of  $R^2=0.84$ , while  
305 the correlation between Beagle predicted alleles and the true calls was  $R^2=0.69$  (Fig. 2  
306 A). At higher depths of coverage (>5X), the raw data provides ample information to  
307 distinguish between homozygous and heterozygous genotypes, allowing Beagle to  
308 determine the correct genotype. The PHG, however, is able to distinguish between the

309 available haplotypes at a coverage of 0.5X and adding additional sequence data does  
310 not increase the accuracy, as there is no correlation between accuracy and coverage  
311 beyond 0.5X.

312 The improved performance of the PHG is most noticeable in its accurate  
313 predictions of heterozygous and rare genotypes. The PHG was able to impute  
314 genotypes with high accuracy regardless of allele class (Figure 2B). The PHG's high  
315 accuracy at low allele frequencies for both homozygous (Figure 2C) and heterozygous  
316 genotypes (Figure 2D), display its ability to impute rare alleles.



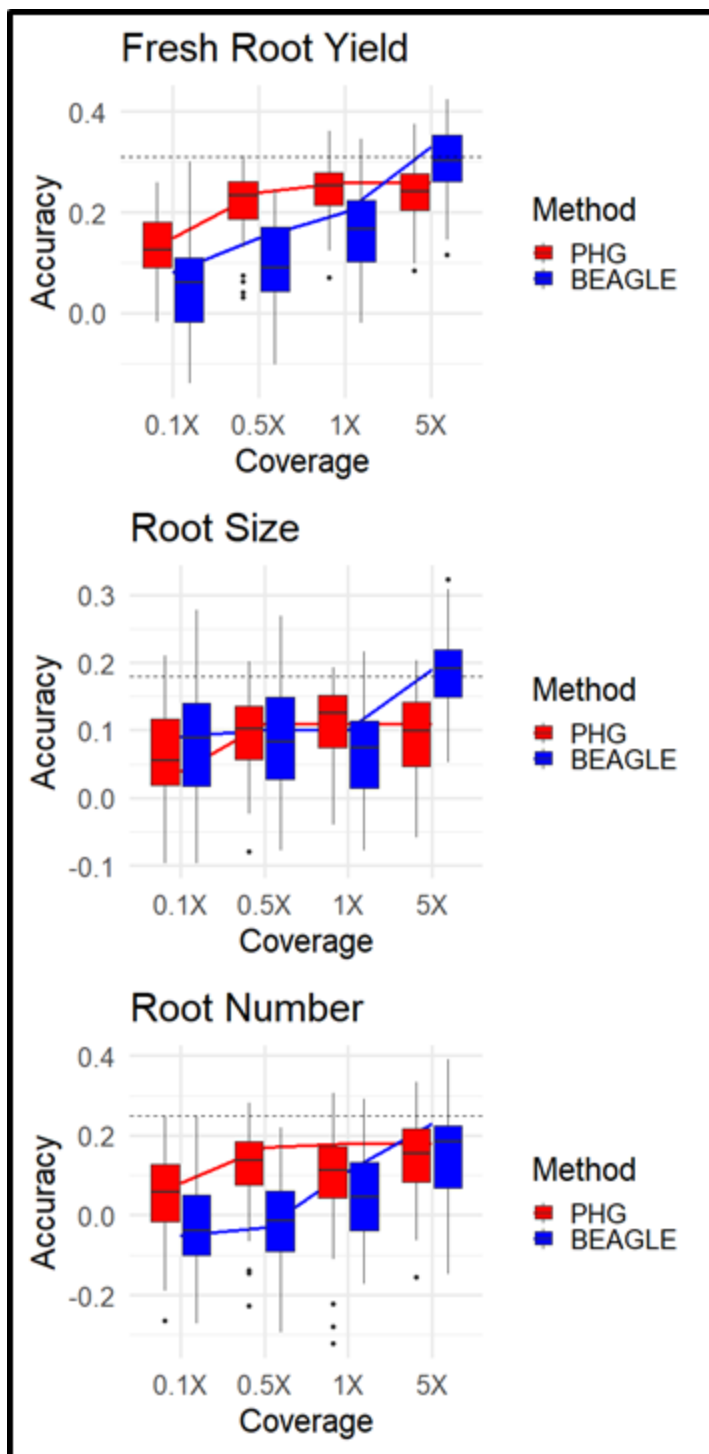
317 **Figure 2. Imputation accuracy from skim sequencing. A) Displays correlation**  
318 **between imputed and true variants by imputing with the PHG and Beagle at**  
319 **different levels of skim sequencing. B) Displays concordance between true and**  
320 **imputed allele at 1X coverage separated by alleles classes: minor, heterozygous,**  
321 **and major. C) Imputation accuracy at 1X coverage is shown for homozygous**  
322 **genotypes separated by allele frequency of the true allele a that locus. D)**



324 **Imputation accuracy at 1X coverage is shown for heterozygous genotypes**  
325 **separated by minor allele frequency at that locus.**

326

327 The imputed genotypes were then utilized in a genomic prediction scheme  
328 consisting of 57 cassava clones (Supplemental Fig. 5) from a single breeding program.  
329 The related nature of the clones ensured an adequate level of heritability to assess  
330 genomic prediction accuracy. Ten-fold cross validations and single holdout validation  
331 showed that imputation accuracy generally appeared to follow the genomic prediction  
332 accuracy, for fresh root yield and root number, while no clear pattern was apparent for  
333 the root size trait (Fig. 3).



334

335 **Figure 3 Genomic Prediction Cross Validation. 10-Fold cross validation (box) and**

336 **single holdout cross validation (line) show genomic prediction accuracies of 3**

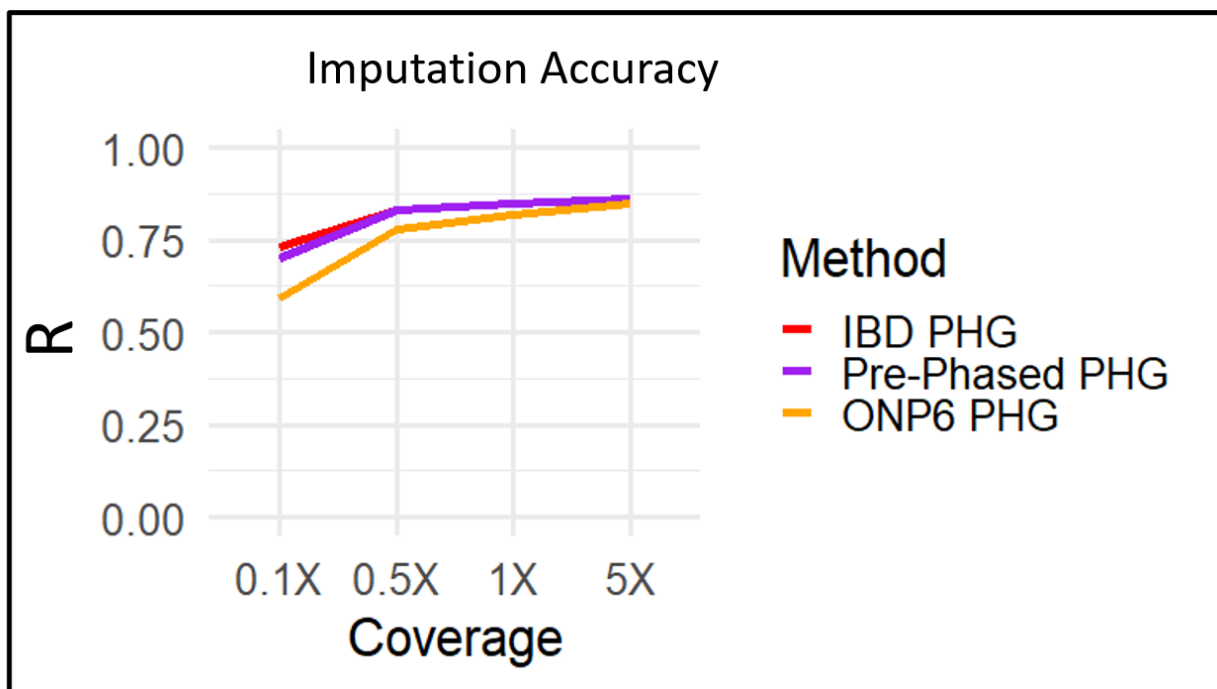
337 **root traits using different imputation methods at varied sequence depths. Single**  
338 **holdout cross validation using complete genotype dataset is shown (dashed line).**

339

### 340 **Phased Haplotype PHG**

341 We tested the viability of populating the PHG with haplotypes phased by other  
342 methods. We compared the IBD method of sampling phased haplotypes to two  
343 methods of phasing variants. The first method used Beagle and HAPCUT2 to phase  
344 the variants called from the HapMapII WGS data. The second method utilized six  
345 cassava clones with ONP long-read data. The IBD method of populating the cassava  
346 PHG produced the highest accuracy (Fig. 4). These results suggest that Beagle and  
347 HAPCUT could not accurately phase heterozygous haplotypes at this scale but  
348 maintained comparable accuracy due to IBD haplotypes. While the PHG was made  
349 from 6 clones with ONP data, it relied on a far narrower set of germplasm. This  
350 suggests that accurate haplotypes are likely captured using this method but lack  
351 adequate sampling to capture sufficient haplotypes.

352

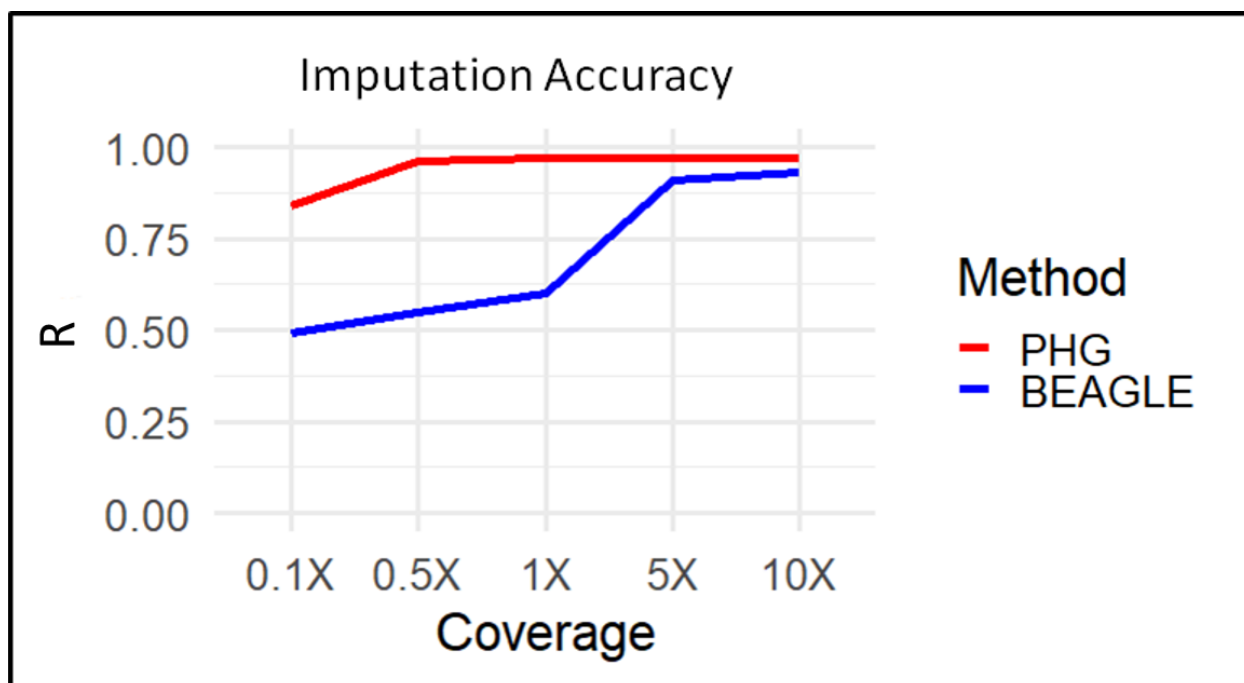


353  
354 **Figure 4. Haplotype phasing methods in the PHG. Imputation accuracy is shown**  
355 **for three different methods of populating a PHG. First the IBD PHG (red) was**  
356 **populated using homozygous haplotypes from the 241 HapMapII clones. Second,**  
357 **the Pre-Phased PHG (Purple) used Beagle and HPACUT2 to phase these same**  
358 **clones. Third, the ONP6 PHG (Yellow) used ONP long-reads and WhatsHap to**  
359 **phase six related taxa to the test set.**

360  
361 **Imputation Simulation**

362 Evident from the tests using haplotypes from IBD regions of the genome,  
363 sampling phased haplotypes is a difficult aspect of creating an effective PHG in a  
364 heterozygous species. To explore the performance of the PHG in a scenario where one  
365 could aptly sample the diversity of haplotypes, we used simulated offspring from a set of  
366 20 phased genomes. While phasing errors exist, we accepted these phases as truth for  
367 the simulation of offspring. This ensured that all haplotypes present in the offspring

368 exist in the PHG database. We found that the disparity in accuracies between PHG and  
369 Beagle at high sequence coverage disappeared in our simulation (Fig. 5), while the  
370 trend in Beagle accuracy was very similar to our empirical tests. While the simulation  
371 does represent an ideal scenario, including a narrower set of germplasm, it highlights  
372 the performance of the PHG when accurately phased haplotypes are available.  
373



374  
375 **Figure 5 Imputation accuracy with simulated genotypes. A simulated scenario**  
376 **where a PHG is populated using 20 parents with full phased information.**  
377 **Correlation between imputed and true variants by imputing with the PHG and**  
378 **Beagle at different levels of skim sequencing.**

379

380

381

382

## DISCUSSION

383           We have detailed a method of implementing a PHG for the heterozygous plant  
384 species cassava. This PHG database utilizes phased haplotypes to predict missing  
385 genotypes from low depth input sequence. Runs of homozygosity formed by IBD  
386 relationships proved to be the most reliable method of sampling phased haplotypes  
387 given the available data (Fig. 4). This method of obtaining haplotypes, while not able  
388 obtain the full diversity of alleles, captures 77% of common alleles and produces ample  
389 haplotypes for significant imputation accuracy at very low sequence depth (Fig. 2A).

390           The high accuracy of the PHG demonstrates its potential as an imputation tool  
391 for use in heterozygous crops. The advantages of the PHG imputation methodology are  
392 especially evident in its accuracy at calling rare and heterozygous alleles (Fig 2C,2D).  
393 Furthermore, the observed weaker relationship between allele frequency and imputation  
394 accuracy, highlights its ability to predict rare alleles. Across both simulated and  
395 empirical experiments, we found that the ability of the PHG to impute whole genome  
396 variants was consistent at or above 0.5X sequence coverage. The haplotype-based  
397 representation of the genome enables this imputation methodology to overcome the  
398 logistical hurdles such as those produced by sequencing and assembly errors, repetitive  
399 sequences, and poor alignments.

400           The plateau reached in imputation accuracy (Fig. 2A) using the PHG most likely  
401 indicates that we have not sufficiently sampled the diversity of possible haplotypes. At  
402 sequence coverages of 5X and more, the raw data can produce the true genotypes and  
403 little imputation of missing genotypes is occurring. The disparity between the PHG and  
404 Beagle at these high coverages shows the presence of missing haplotypes in the  
405 database, rather than any disparity in performance. This is supported by a visible

406 relationship between homozygous incidence in our population and reference range  
407 imputation accuracy (Supplemental Fig. 6), suggesting that those ranges with poor  
408 imputation accuracy were not amply sampled. The length and abundance of the IBD  
409 runs of homozygosity in our dataset likely determine the ability of the HMM to accurately  
410 predict haplotypes. This hypothesis for decreased accuracy is also supported by the  
411 removal of such a disparity under simulation, where all possible haplotypes are sampled  
412 in the database (Fig. 5). These results suggest that, although an already powerful tool,  
413 the PHG achieves maximum performance with sufficient sampling of available  
414 haplotypes.

415 While the imputation accuracy of the PHG is limited based on the haplotype  
416 sampling, its high accuracy with low levels of input sequence highlights its potential for  
417 genomic applications, where sparse genotyping is common. We showed that this is  
418 true regarding genomic prediction by performing cross-validations with the imputed  
419 genotypes (Fig. 3). The genomic prediction was still limited by imputation accuracy, but  
420 by enabling higher accuracy we can achieve more reliable predictions (Pimentel *et al.*  
421 2015; Wang *et al.* 2016; Van Den Berg *et al.* 2017).

422 With increased environmental pressures in a changing global climate and  
423 growing populations, accelerated breeding is vital to sustainable food production. With  
424 increased imputation accuracy from more limited genotyping inputs, smaller breeding  
425 programs can afford to implement GS, enabling them to increase selection pressure  
426 across their breeding pools. Similarly, imputation to genome-wide scale can bridge gaps  
427 between different data sets containing information on different marker panels, enabling  
428 the use of larger datasets for prediction. Accurate imputation could also enable





452

## ACKNOWLEDGMENTS

453           This study is made possible by the funding and support of the Nextgen Cassava  
454           project, the Bill and Malinda Gates foundation, and the USDA-ARS. We also  
455           acknowledge the programming staff in the Buckler lab who created and support the  
456           development of the practical haplotype graph. Lastly, we are grateful for the greater  
457           Nextgen Cassava community for supporting the curation of genotype and phenotype  
458           data used in this project as well as the organization of this data in [cassavabase.org](http://cassavabase.org).

459

460 REFERENCES

- 461 Alipour, H., G. Bai, G. Zhang, M. R. Bihamta, V. Mohammadi *et al.*, 2019 Imputation  
462 accuracy of wheat genotyping-by-sequencing (GBS) data using barley and wheat  
463 genome references. *PLoS One* 14:.
- 464 Van Den Berg, I., P. J. Bowman, I. M. MacLeod, B. J. Hayes, T. Wang *et al.*, 2017 Multi-  
465 breed genomic prediction using Bayes R with sequence data and dropping variants  
466 with a small effect. *Genet. Sel. Evol.* 49: 1–15.
- 467 Browning, B. L., Y. Zhou, and S. R. Browning, 2018 A One-Penny Imputed Genome  
468 from Next-Generation Reference Panels. *Am. J. Hum. Genet.* 103: 338–348.
- 469 Cleveland, M. A., J. M. Hickey, and B. P. Kinghorn, 2011 Genotype imputation for the  
470 prediction of genomic breeding values in non-genotyped and low-density  
471 genotyped individuals, pp. S6 in *BMC Proceedings*, BioMed Central.
- 472 Edge, P., V. Bafna, and V. Bansal, 2017 HapCUT2: Robust and accurate haplotype  
473 assembly for diverse sequencing technologies. *Genome Res.* 27: 801–812.
- 474 Fang, L., G. Sahana, P. Ma, G. Su, Y. Yu *et al.*, 2017 Exploring the genetic architecture  
475 and improving genomic prediction accuracy for mastitis and milk production traits in  
476 dairy cattle by mapping variants to hepatic transcriptomic regions responsive to  
477 intra-mammary infection. *Genet Sel Evol* 49: 44.
- 478 Fragoso, C. A., C. Heffelfinger, H. Zhao, and S. L. Dellaporta, 2016 Imputing Genotypes  
479 in Biallelic Populations from Low-Coverage Sequence Data. *Genetics* 202: 487–  
480 495.
- 481 Franco, J. A. V., J. L. Gage, P. J. Bradbury, L. C. Johnson, Z. R. Miller *et al.*, 2020 A  
482 Maize Practical Haplotype Graph Leverages Diverse NAM Assemblies. *bioRxiv*

483           2020.08.31.268425.

484   de Freitas, J. P. X., V. da Silva Santos, and E. J. de Oliveira, 2016 Inbreeding  
485           depression in cassava for productive traits. *Euphytica* 209: 137–145.

486   Friedenberg, S. G., and K. M. Meurs, 2016 Genotype imputation in the domestic dog.  
487           *Mamm. Genome* 27: 485–494.

488   Gilmour, A. R., 2019 Average information residual maximum likelihood in practice. *J.*  
489           *Anim. Breed. Genet.* 136: 262–272.

490   Haller, B. C., and P. W. Messer, 2019 Evolutionary Modeling in SLiM 3 for Beginners.  
491           *Mol. Biol. Evol.* 36: 1101–1109.

492   Heffner, E. L., M. E. Sorrells, and J. L. Jannink, 2009 Genomic selection for crop  
493           improvement. *Crop Sci.* 49: 1–12.

494   Jensen, S. E., J. R. Charles, K. Muleta, P. J. Bradbury, T. Casstevens *et al.*, 2020 A  
495           sorghum practical haplotype graph facilitates genome-wide imputation and cost-  
496           effective genomic prediction. *Plant Genome* 1–15.

497   Kendig, K. I., S. Baheti, M. A. Bockol, T. M. Drucker, S. N. Hart *et al.*, 2019 Sentieon  
498           DNaseq Variant Calling Workflow Demonstrates Strong Computational  
499           Performance and Accuracy. *Front. Genet.* 10: 736.

500   Kono, T. J. Y., C. Liu, E. E. Vonderharr, D. Koenig, J. C. Fay *et al.*, 2019 The Fate of  
501           Deleterious Variants in a Barley Genomic Prediction Population. *Genetics* 213:  
502           1531–1544.

503   Kremling, K. A. G., S. Y. Chen, M. H. Su, N. K. Lepak, M. C. Romay *et al.*, 2018  
504           Dysregulation of expression correlates with rare-allele burden and fitness loss in  
505           maize. *Nature* 555: 520–523.

- 506 Loh, P. R., P. F. Palamara, and A. L. Price, 2016 Fast and accurate long-range phasing  
507 in a UK Biobank cohort. *Nat. Genet.* 48: 811–816.
- 508 MacLeod, I. M., P. J. Bowman, C. J. Vander Jagt, M. Haile-Mariam, K. E. Kemper *et al.*,  
509 2016 Exploiting biological priors and sequence variants enhances QTL discovery  
510 and genomic prediction of complex traits. *BMC Genomics* 17: 1–21.
- 511 Meuwissen, T. H. E., B. J. Hayes, and M. E. Goddard, 2001 Prediction of total genetic  
512 value using genome-wide dense marker maps. *Genetics* 157: 1819–1829.
- 513 Nazzicari, N., F. Biscarini, P. Cozzi, E. C. Brummer, and P. Annicchiarico, 2016 Marker  
514 imputation efficiency for genotyping-by-sequencing data in rice (*Oryza sativa*) and  
515 alfalfa (*Medicago sativa*). *Mol. Breed.* 36: 69.
- 516 Parmar, A., B. Sturm, and O. Hensel, 2017 Crops that feed the world: Production and  
517 improvement of cassava for food, feed, and industrial uses. *Food Secur.* 9: 907–  
518 927.
- 519 Pimentel, E. C. G., C. Edel, R. Emmerling, and K. U. Götz, 2015 How imputation errors  
520 bias genomic predictions. *J. Dairy Sci.* 98: 4131–4138.
- 521 Ramu, P., W. Esuma, R. Kawuki, I. Y. Rabbi, C. Egesi *et al.*, 2017 Cassava haplotype  
522 map highlights fixation of deleterious mutations during clonal propagation. *Nat.*  
523 *Genet.* 49: 959–963.
- 524 Romay, M. C., 2018 Rapid, Affordable, and Scalable Genotyping for Germplasm  
525 Exploration in Maize, pp. 31–46 in Springer, Cham.
- 526 Schrunner, S. D., R. S. Mari, J. Ebler, M. Rautiainen, L. Seillier *et al.*, 2020 Haplotype  
527 Threading: Accurate Polyploid Phasing from Long Reads. *bioRxiv*  
528 2020.02.04.933523.

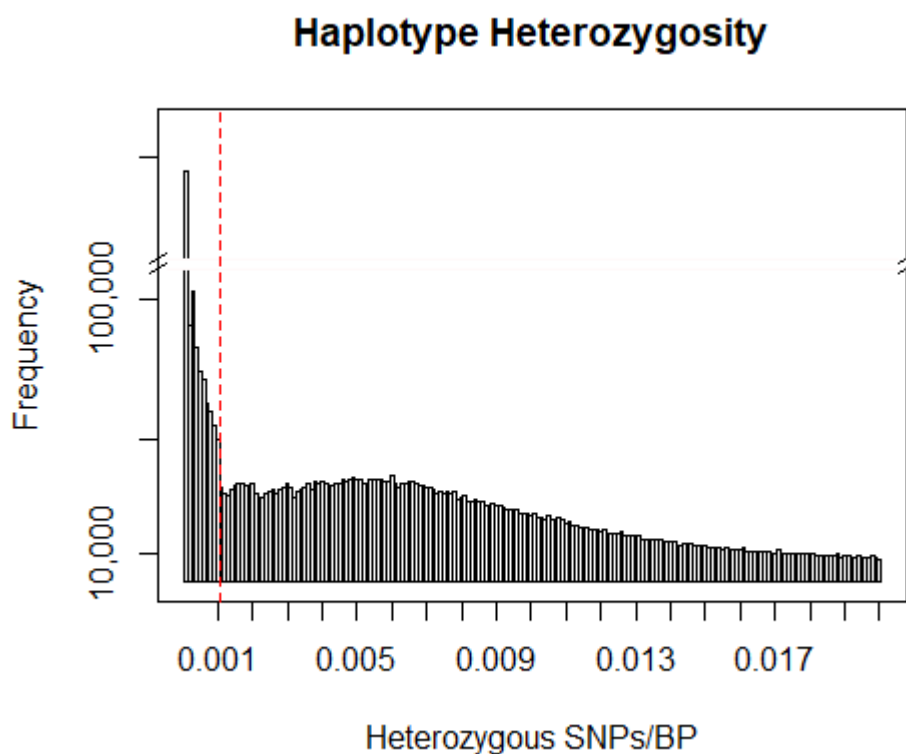
- 529 Spencer, C. C. A., Z. Su, P. Donnelly, and J. Marchini, 2009 Designing Genome-Wide  
530 Association Studies: Sample Size, Power, Imputation, and the Choice of  
531 Genotyping Chip (J. D. Storey, Ed.). PLoS Genet. 5: e1000477.
- 532 Stephens, Z. D., M. E. Hudson, L. S. Mainzer, M. Taschuk, M. R. Weber *et al.*, 2016  
533 Simulating next-generation sequencing datasets from empirical mutation and  
534 sequencing models. PLoS One 11:.
- 535 Swarts, K., H. Li, J. A. Romero Navarro, D. An, M. C. Romay *et al.*, 2015 Novel  
536 Methods to Optimize Genotypic Imputation for Low-Coverage, Next-Generation  
537 Sequence Data in Crop Plants. Plant Genome 7: 0.
- 538 Torkamaneh, D., B. Boyle, and F. Belzile, 2018 Efficient genome-wide genotyping  
539 strategies and data integration in crop plants. Theor. Appl. Genet. 131: 499–511.
- 540 Wang, Y., G. Lin, C. Li, and P. Stothard, 2016 Genotype Imputation Methods and Their  
541 Effects on Genomic Predictions in Cattle. Springer Sci. Rev. 4: 79–98.
- 542 Xu, Y., X. Liu, J. Fu, H. Wang, J. Wang *et al.*, 2020 Enhancing Genetic Gain through  
543 Genomic Selection: From Livestock to Plants. Plant Commun. 1: 100005.
- 544 Yang, J., S. Mezmouk, A. Baumgarten, E. S. Buckler, K. E. Guill *et al.*, 2017 Incomplete  
545 dominance of deleterious alleles contributes substantially to trait variation and  
546 heterosis in maize (J. C. Fay, Ed.). PLOS Genet. 13: e1007019.
- 547 Yun, L., C. Willer, S. Sanna, and G. Abecasis, 2009 Genotype imputation. Annu. Rev.  
548 Genomics Hum. Genet. 10: 387–406.
- 549 Zheng, X., D. Levine, J. Shen, S. M. Gogarten, C. Laurie *et al.*, 2012 A high-  
550 performance computing toolset for relatedness and principal component analysis of  
551 SNP data. Bioinformatics 28: 3326–3328.

552

553

554

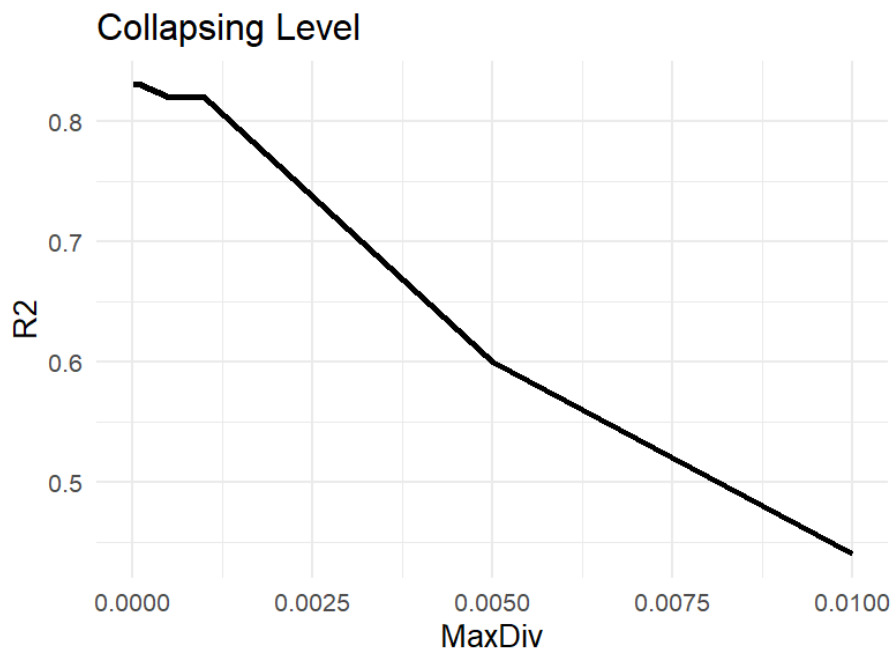
## SUPPLEMENTAL FIGURES



555

---

556 **Supplemental Figure 1. Histogram of sampled haplotypes by the number of**  
557 **heterozygous SNPs per base pair. Dotted line shows the threshold chosen to**  
558 **distinguish nearly IBD haplotypes.**

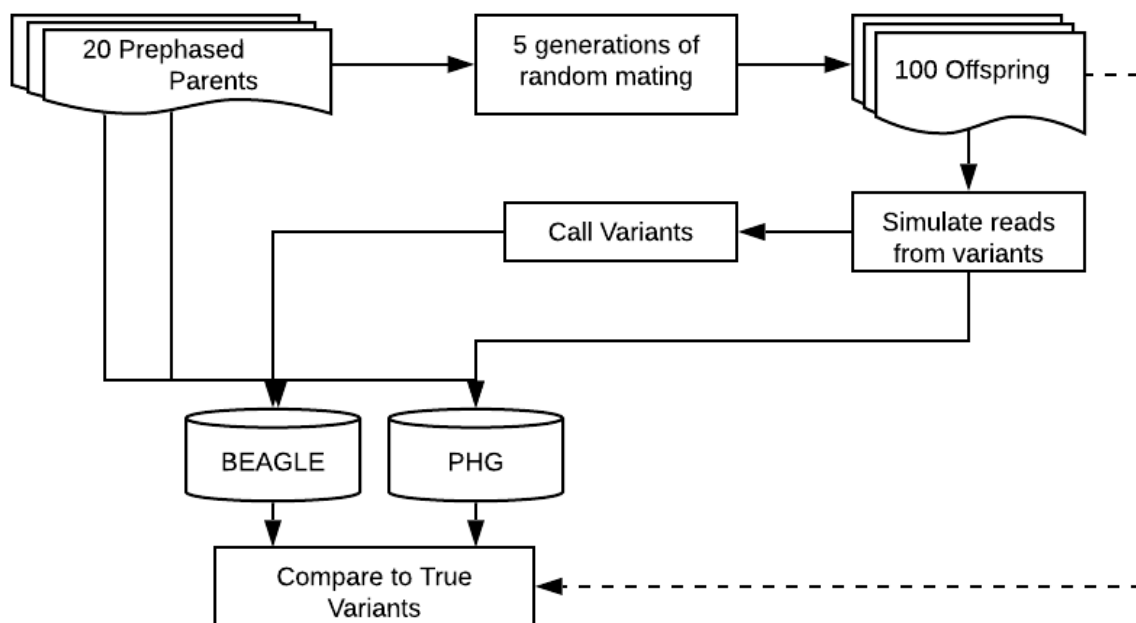


559

560 **Supplemental Figure 2. Correlations of Imputed calls to true calls at different**  
561 **haplotype collapsing levels based on the maximum divergence parameter of the**  
562 **PHG.**

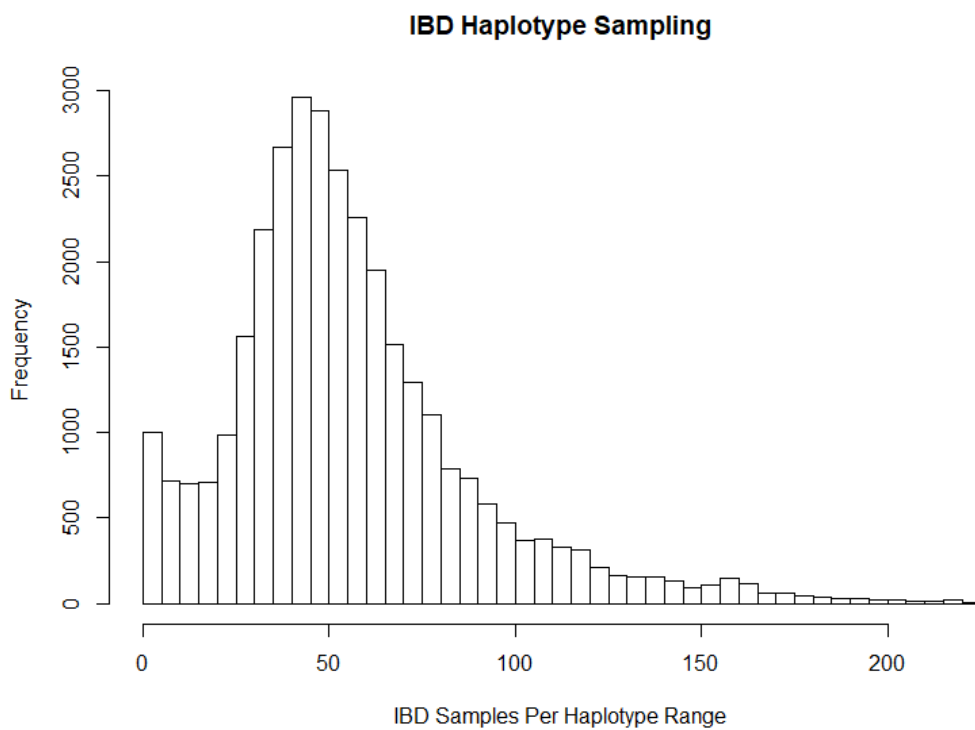
563





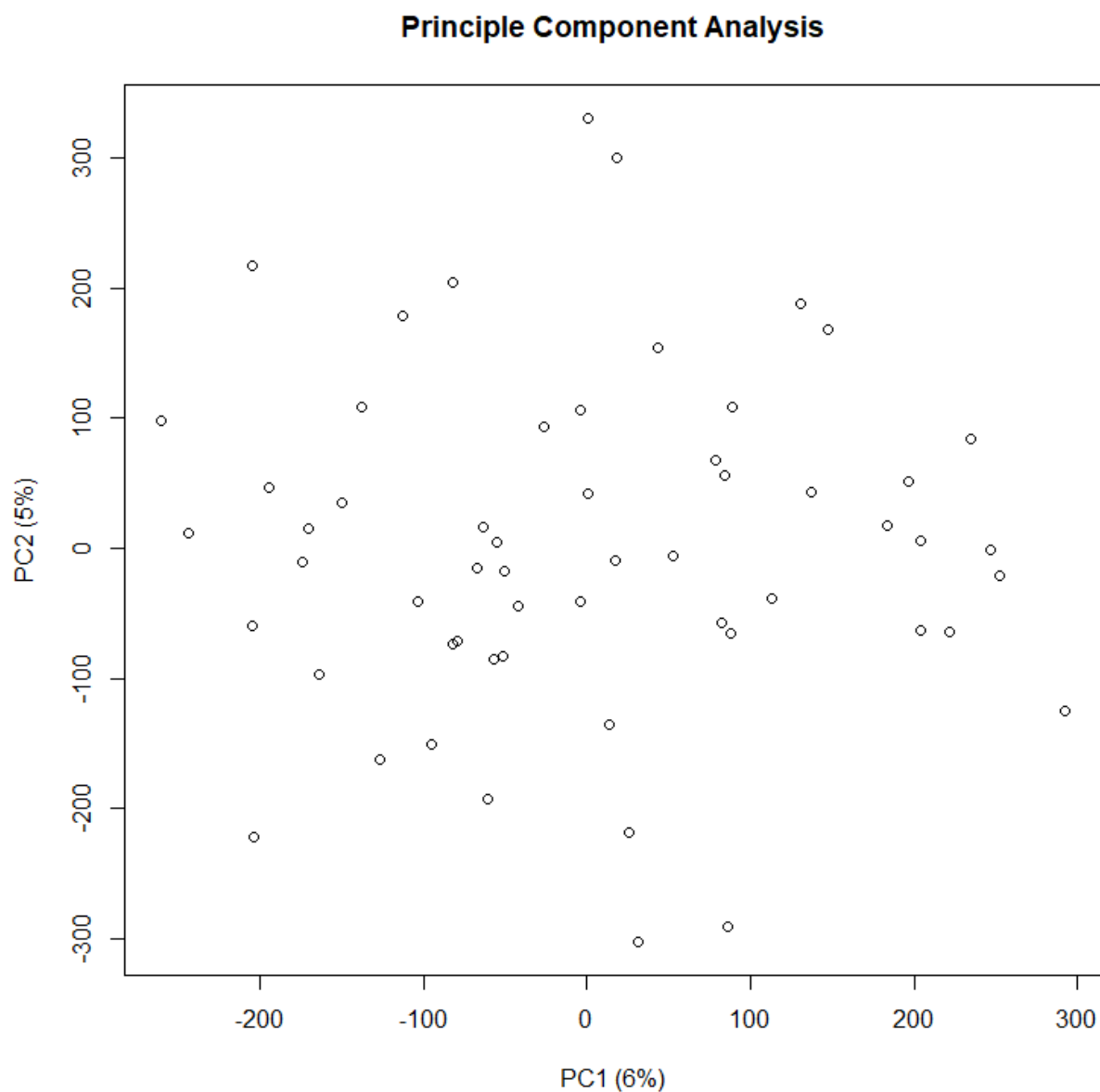
564

565 **Supplemental Figure 3. Simulation and Imputation Schema. Phased genomes**  
566 **from 20 parents from single breeding program were used to simulate generations**  
567 **of mating and recombination.**



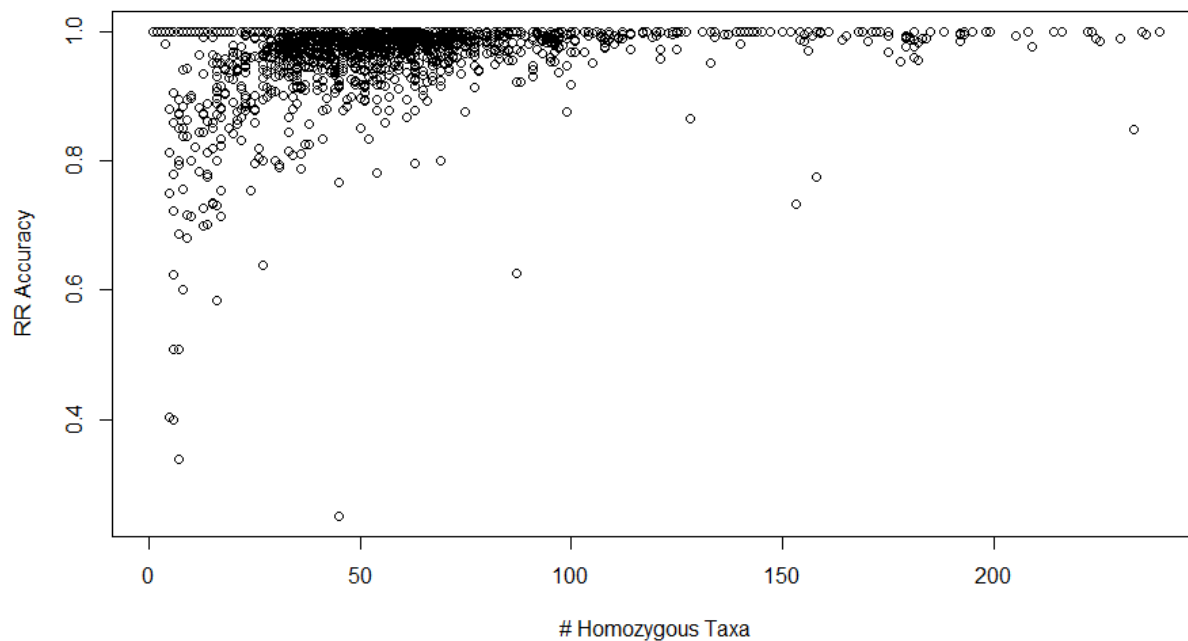
568

569 **Supplemental Figure 4. Histogram of IBD sampling frequency of all reference**  
570 **ranges. Y-axis shows the number of reference ranges with a given number of IBD**  
571 **samples.**



572

573 **Supplemental Figure 5. Principal component analysis of 57 clones used in**  
574 **genomic prediction cross validation. Lack of clusters show little population**  
575 **structure among the clones.**



576

577 **Supplemental Figure 6. Imputation accuracy at each reference (y-axis) range by**  
578 **homozygous incidence in the HapMapII Population (x-axis). Low accuracy shown**  
579 **at reference ranges with low incidence of homozygous taxa.**

580

581