

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, USA

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

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

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

ABSTRACT

20
21 Genomic applications such as genomic selection and genome-wide association
22 have become increasingly common since the advent of genome sequencing. The cost of
23 sequencing has decreased in the past two decades, however genotyping costs are still
24 prohibitive to gathering large datasets for these genomic applications, especially in non-model
25 species where resources are less abundant. Genotype imputation makes it possible to infer
26 whole genome information from limited input data, making large sampling for genomic
27 applications more feasible. Imputation becomes increasingly difficult in heterozygous species
28 where haplotypes must be phased. The Practical Haplotype Graph is a recently developed tool
29 that can accurately impute genotypes, using a reference panel of haplotypes. We showcase
30 the ability of the Practical Haplotype Graph to impute genomic information in the highly
31 heterozygous crop cassava (*Manihot esculenta*). Accurately phased haplotypes were sampled
32 from runs of homozygosity across a diverse panel of individuals to populate PHG, which proved
33 more accurate than relying on computational phasing methods. The Practical Haplotype Graph
34 achieved high imputation accuracy, using sparse skim-sequencing input, which translated to
35 substantial genomic prediction accuracy in cross validation testing. The Practical Haplotype
36 Graph showed improved imputation accuracy, compared to a standard imputation tool Beagle,
37 especially in predicting rare alleles.

INTRODUCTION

38
39 The past decade has seen an abundance of genomic sequence data produced
40 for research and application in agricultural crops. With these new technologies, comes
41 questions on how to effectively implement them (Torkamaneh *et al.* 2018). Two of the
42 most common uses of genome-wide sequence data are genomic selection (GS) and
43 genome-wide association studies (GWA). While most GWAS attempt to locate distinct,
44 causative regions of the genome, genomic selection incorporates all available markers

45 to predict plant traits (Meuwissen *et al.* 2001). Genomic selection leverages a training
46 set population that has both genotypic and phenotypic data to predict traits in a related
47 germplasm with only genotypic data (Heffner *et al.* 2009). This allows breeders to both
48 increase accuracy in selecting traits with low heritability and accelerate the rate of
49 selections by decreasing selection cycle time (Xu *et al.* 2020).

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

63 To affordably obtain high density genotypes or to bridge information between
64 different marker platforms it becomes necessary to impute missing genotypes from
65 available genotype data. Increasing the stability across genotyping platforms and
66 reducing per-sample costs, becomes even more relevant in plant breeding scenarios,
67 where many thousands of offspring are evaluated and changes in marker platform are

68 common. Computational techniques to impute genome-wide information have been
69 produced to bridge genotypic information from different marker panels and augment
70 genotypic information from limited inputs (Yun *et al.* 2009). Genomic imputation
71 methods often rely on a related training set with high confidence genotypic information
72 to then predict missing genotypes. These methods have been shown to improve
73 consistency and efficiency of analyses of both genome wide associations (Spencer *et*
74 *al.* 2009) and genomic selection (Cleveland *et al.* 2011).

75 Imputation is very common in genomic studies but is still plagued with barriers to
76 high accuracy in many species. Known limitations of imputation stem from LD, allele
77 frequencies, and population structure of the training population (Alipour *et al.* 2019).
78 These difficulties are further compounded when working with a highly heterozygous
79 crop, where both copies of the genome need to be modeled (Fragoso *et al.* 2016;
80 Nazzicari *et al.* 2016). Heterozygosity introduces the challenge of phasing, the process
81 assigning alleles to haplotypes, a challenge that is not limited to plants (Friedenberg
82 and Meurs 2016). Imputation accuracy has been shown to affect the accuracy of
83 genomic prediction in multiple scenarios (Pimentel *et al.* 2015; Wang *et al.* 2016; Van
84 Den Berg *et al.* 2017). Additionally, when tracking causative variation through the
85 genome, high accuracy in imputation is necessary to evaluate variation across the
86 entire genome. Highly accurate imputation methods are needed to increase the gains
87 made by genomic selection by making genotyping cheaper, more accurate, and more
88 consistent.

89 It has been shown that rare variants contribute to the genetic load and overall
90 performance of crops (Yang *et al.* 2017; Kremling *et al.* 2018; Kono *et al.* 2019), making

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

102 In maize, it's been shown that Beagle has difficulty accurately imputing rare
103 variants, while a haplotype library based method such as FILLIN can do so more easily
104 (Swarts *et al.* 2015). A recently developed method known as the Practical Haplotype
105 Graph (PHG) was created to leverage known haplotypes in a graph structure to
106 efficiently impute genotypes. The PHG simplifies the genome to a set of distinct regions
107 of the genome, for which it defines haplotypes. These haplotypes are constructed from
108 whole genome sequence data or genome assemblies and are used to construct a trellis
109 graph, capturing the diversity of haplotypes at each range and the relationships
110 between adjacent haplotype regions. Sequence reads are then aligned to the graph
111 and an HMM is applied to predict the most likely haplotypes. By aligning reads to pan-
112 genome haplotypes, the PHG minimizes errors due to reference bias, poor alignment,
113 and mis-called variants. Utilizing a PHG methodology in plant and animal applications

114 can improve the quality and quantity of genotype data for use in breeding and mapping
115 scenarios.

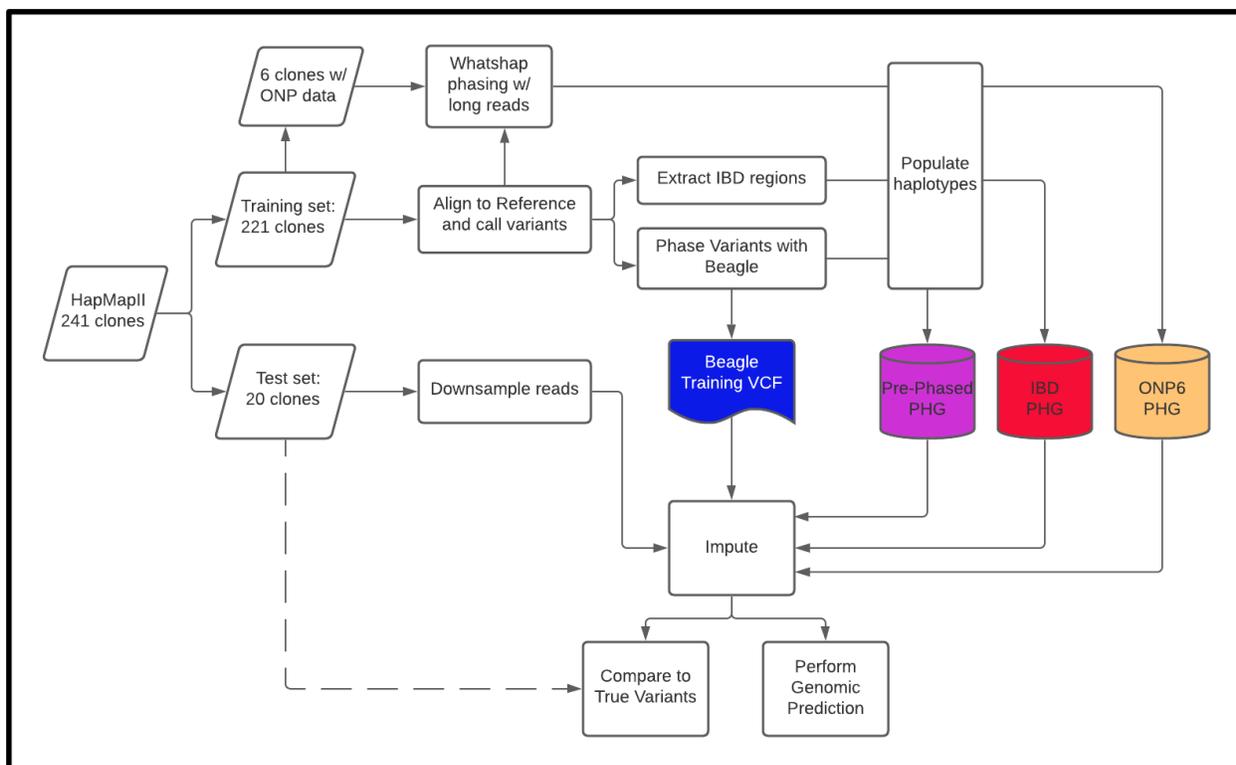
116 Here we showcase the potential application of the PHG in imputation of
117 heterozygous crops. The PHG has already been shown to be an efficient tool for aiding
118 imputation and genomic selection in breeding of the inbred cereal crop Sorghum
119 (Jensen *et al.* 2020). It has also been implemented to impute genotypes in highly
120 diverse maize lines (Valdes Franco 2020). To show the utility of the PHG in a
121 heterozygous crop we must overcome two distinct challenges: obtaining phased
122 haplotypes to populate the database and 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. In this study we
130 utilize sequence data from the previously published HapMapII in cassava (Ramu *et al.*
131 2017), which includes WGS data for 241 cassava clones. This data is used to produce
132 a PHG in cassava and showcase its effectiveness in genomic imputation in a
133 heterozygous crop. We further validate these methods through genomic prediction and
134 simulation.

135

136

MATERIALS AND METHODS



137

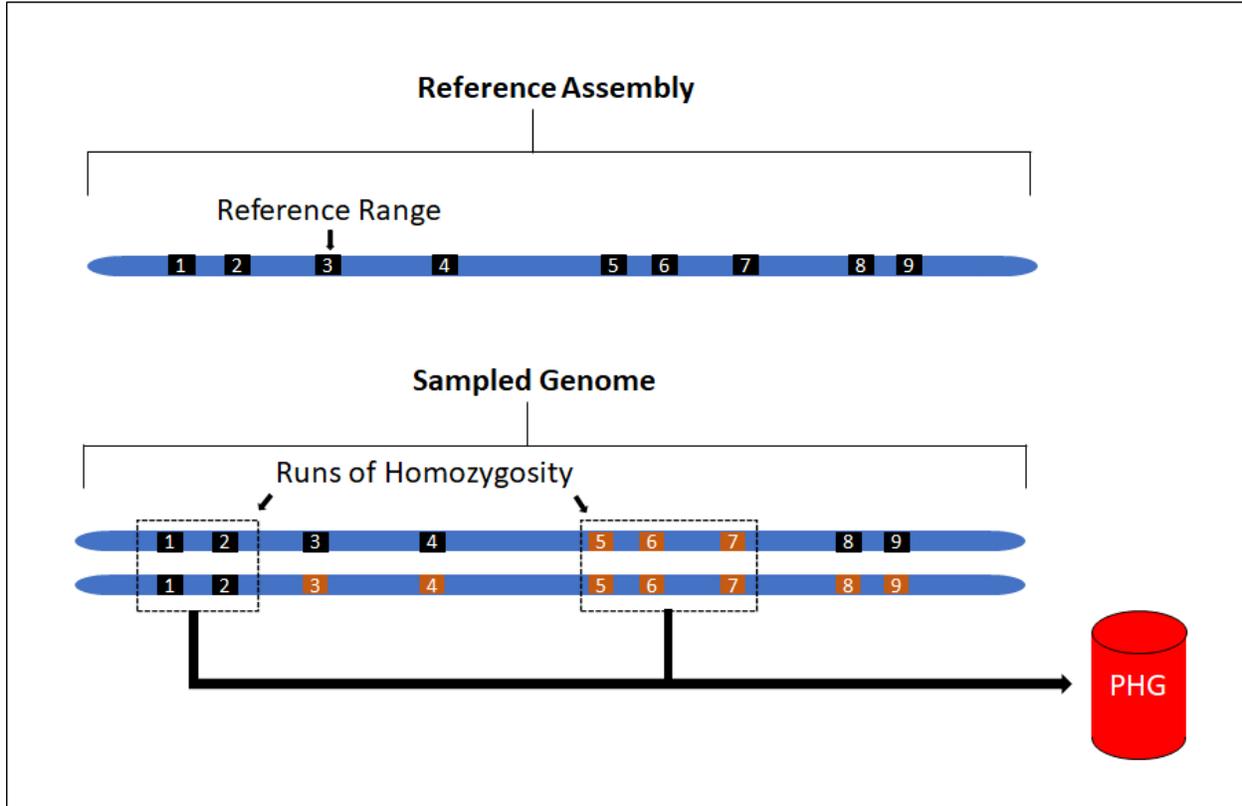
138 **Figure 1. Imputation Methodology Flowchart. Diagram of methods used for the**
 139 **building PHG databases and performing imputation evaluations.**

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. We used the cassava v6 reference genome assembly in this study, which
 147 contains 18 chromosome level scaffolds summing to ~500Mbp of the estimated genome
 148 size of 700Mbp. Haplotype regions, termed here as reference ranges, were defined by
 149 genic regions with additional 1000bp flanking sequence resulting in ~32,000 reference
 150 ranges after merging overlapping ranges, with an average size of 4kbp.

151 The detailed process of creating a PHG is outlined at
152 ["https://bitbucket.org/bucklerlab/practicalhaplotypegraph/wiki/Home"](https://bitbucket.org/bucklerlab/practicalhaplotypegraph/wiki/Home) and has been
153 described previously (Jensen *et al.* 2020; Valdes Franco 2020). Here, we outline the
154 specific steps taken to create a PHG in the heterozygous crop cassava (Fig. 1). The
155 major hurdle to producing a haplotype graph in a heterozygous species is obtaining
156 accurately phased haplotypes. Because many of these cassava lines are cultivated
157 taxa, we expected to find identical by descent (IBD) haplotypes brought about by
158 generations of breeding within restricted breeding pools. These IBD segments provide
159 confidently phased haplotypes as well as capturing their relationships to adjacent
160 haplotypes (Fig. 2). We identified and sampled these homozygous haplotypes which
161 we inferred to represent IBD haplotypes. This was done by measuring the number of
162 heterozygous variants for each reference range in each taxon, then classifying those
163 haplotypes as homozygous or not. The threshold for haplotypes to be considered IBD
164 was determined empirically to be 0.001 heterozygous SNPs per base pair
165 (Supplemental Fig. 1), as *de novo* mutations or errors in variant calling may produce low
166 levels of perceived heterozygosity. This threshold was additionally validated by testing

167 imputation accuracy of the PHG.



168

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

174

175 After haplotypes were sampled from IBD regions of the genome, they were
176 loaded as GVCF files into a PHG database. Similar haplotypes were then collapsed
177 based on sequence similarity to produce a representative set of available haplotypes.
178 Haplotypes are collapsed to make alignment more efficient, while retaining as much
179 distinct haplotype information as possible. Collapsing is performed using an

180 unweighted pair group method with arithmetic mean (upgma) tree from pairwise
181 distance matrix from sequence variants to measure the similarity between haplotypes.
182 Based on imputation accuracy tests, we chose a level of similarity (PHG parameter:
183 maximum divergence) to collapse haplotypes of 0.001, corresponding to less than 1 in
184 1000 nucleotide differences between haplotypes. This level of collapsing maintains
185 high accuracy while collapsing redundant haplotypes (Supplemental Fig. 2). We then
186 produced a pan-genome composed of consensus haplotypes representing the diversity
187 of haplotypes.

188 **Predicting Haplotypes**

189 Once we obtained a set of consensus haplotypes, we implemented an HMM to
190 infer genome-wide haplotypes from low depth genotyping data. Sparse genotype
191 information was created by downsampling whole genome sequence data randomly
192 using samtools to simulate skim sequencing. We randomly sampled 20 taxa from the
193 cultivated varieties within the population to serve as a test set for downstream analyses,
194 while using the remaining 221 clones for haplotype sampling. To test different levels of
195 sequencing depth, we down-sampled reads to amounts estimated to represent 0.1X,
196 0.5X, 1X, 5X, and 10X single-end, whole genome sequence coverage. Additionally, we
197 tested imputation using available Genotype-By-Sequencing (GBS) data for these lines.

198 These sampled sequences were aligned to the consensus haplotypes stored in
199 the PHG to impute whole genome variants. A trellis graph is formed with every
200 reference range representing separate ranges and the consensus haplotypes as nodes
201 at each of those ranges. The most likely paths through the graph were then determined
202 using an HMM Viterbi algorithm. Because cassava is heterozygous and diploid, this

203 step produces the two most likely paths for each taxon. The emission and transition
204 probability parameters of the HMM are defined by the genomes of the reference
205 population used to build the database. The emission probabilities are calculated by
206 considering the probability of two given haplotypes, given the aligned reads. The
207 transition probabilities are defined by the edges between haplotypes in the PHG.

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

214 **Beagle imputation**

215 We compared our imputation accuracy results to the common genotype
216 imputation tool Beagle (Browning *et al.* 2018). Beagle was developed for the purpose of
217 human data, but is a common tool used by many plant studies to impute missing
218 genotypes. Because Beagle v4 has the ability to incorporate genotype likelihoods
219 based on read depth, we used it for the imputation of the low depth sequence when it
220 improved accuracy, otherwise we utilized Beagle v5. We used the same HapMapII data
221 from the 241 clones to impute missing genotypes with Beagle.

222 **Genomic Prediction**

223 We used 57 clones from a single breeding program, to reduce effects of
224 population structure, to determine the impact of imputation errors on genomic prediction
225 accuracy using cross validation. Reads were downsampled and imputed as previously

226 described. Three root traits were used for genomic cross validation: fresh root yield,
 227 root size, and root number. Phenotypes for each clone were downloaded from
 228 CassavaBase.org, constituting 57 clones, spanning 23 years from 1996 to 2018, across
 229 13 locations in Africa. Ten-fold cross validation was performed by randomly selecting
 230 10% of the clones to hold out and predict using the remaining clones as a training set.
 231 The correlation between predicted phenotype and the observed best linear unbiased
 232 estimate (BLUE) was used as the prediction accuracy. We performed 50 replications as
 233 well as a single holdout prediction to measure genomic prediction accuracy. A single
 234 step model was performed:

$$235 \quad \hat{y} = \mu + G_i + B_j + R_k + L_l + Y_m + G_iXL_l + G_iXY_m$$

$$236 \quad G_i \sim N(0, G\sigma_G^2), B_j \sim N(0, I\sigma_B^2), R_k \sim N(0, I\sigma_R^2), L_l \sim N(0, I\sigma_L^2), Y_m \sim N(0, I\sigma_m^2)$$

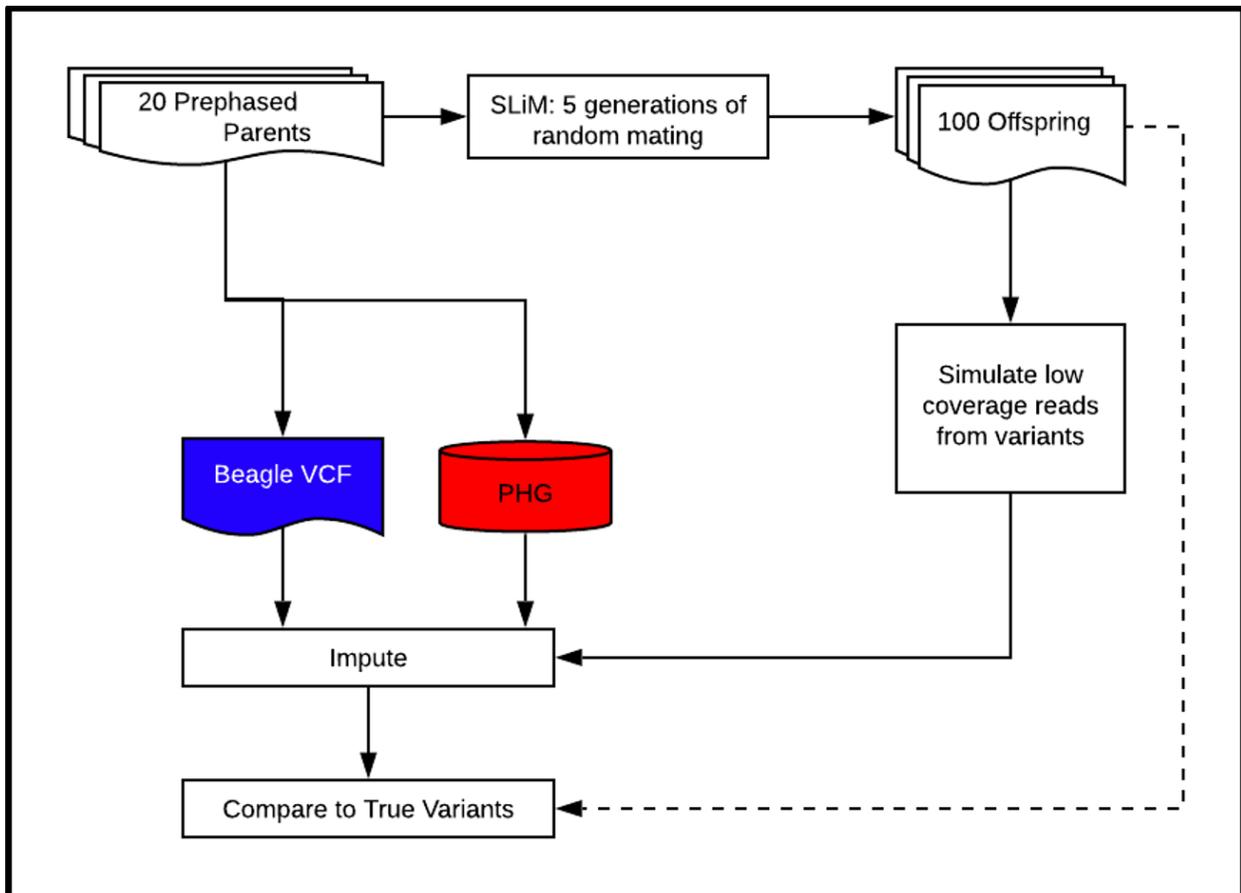
237 Here, \hat{y} is the predicted trait and μ is the fixed effect of the overall mean. Random
 238 effects were fitted as follows: \mathbf{G} is genotype effect of the i th clone, \mathbf{B} is the effect of the
 239 j th block, \mathbf{R} is the effect of the k th replicate, \mathbf{L} is the location of the l th location, \mathbf{Y} is the
 240 effect of the m th year, \mathbf{GXL} is the interactive effect of the i th clone and the l th location,
 241 and \mathbf{GXY} is the interaction effect of the i th clone and the m th year. This was performed
 242 using the mixed model tool ECHIDNA (Gilmour 2019).

243 **Pre-phased Haplotype PHG**

244 We investigated the viability of using computationally phased haplotypes to
 245 curate a PHG database rather than relying on IBD regions of the genome. First we
 246 phased the variants from the 241 cassava clones using a combination of Beagle
 247 (Browning *et al.* 2018) and HAPCUT2 (Edge *et al.* 2017). These variants were used to
 248 create a PHG to be tested against the IBD version of the PHG. The second test utilized

249 Oxford Nanopore (ONP) long-read sequencing from 6 cassava clones within the HMII
250 population. High molecular weight DNA was extracted from young cassava leaves,
251 selected for fragments 20-80 kbp long, and sequenced with MinION following the
252 manufacturer recommendations. Variants were called using Guppy and their variants
253 phased with WhatsHap (Schrinner *et al.* 2020). These 6 clones were then used to
254 populate another PHG, we will identify as the “ONP6 PHG”. Larger reference ranges
255 were divided into smaller regions to increase the probability of sampling correctly
256 phased haplotypes. Twenty clones with the highest relationship to the 6 taxa with ONP
257 data were used as the test set for these tests.

258 Imputation from Simulated Genotypes



259

260 **Figure 3. Simulation Methodology Flowchart. Diagram of simulation scheme**
261 **showing how simulated offspring were generated and used to test imputation**
262 **accuracy under ideal haplotype sampling scenarios.**

263 A sample of 20 related individuals from the HapMapII population were selected
264 to serve as parents for a simulated genotyping scenario. The genomes were phased
265 using Beagle and then used to populate a PHG database. We then used these parents
266 to simulate 5 generations of random mating given a population size of 100 (Fig. 3).
267 Forward genetic simulations were completed using SLiM (Haller and Messer 2019).
268 Artificial short read-sequencing was then simulated for these offspring using neat-
269 genreads (Stephens *et al.* 2016) at varied coverage levels. Reads were then aligned
270 using bwa used to call and impute variants using Sentieon (Kendig *et al.* 2019) and
271 Beagle. Reads were also aligned to the PHG formed from the original parents for
272 imputation.

273 RESULTS

274 Haplotype Sampling

275 To obtain phased haplotypes for the PHG we sampled haplotypes from
276 homozygous regions of each clone. Centuries of clonal propagation and reported
277 inbreeding depression (de Freitas *et al.* 2016) suggest cassava germplasm would be
278 highly heterozygous, however, we found that, on average, ~20% of all reference ranges
279 from each taxon were homozygous. This resulted in a high number of missing
280 haplotypes in each taxon, but a high confidence in the phased haplotypes that were
281 sampled. Despite the variability in the number of homozygous samples by reference
282 range, >90% of the reference ranges were homozygous in at least 10% of the HapMapII

283 population (Supplemental Fig. 3). From these IBD haplotypes we were able to sample
284 ~50% of the segregating sites. This proportion increased to 77% when considering
285 sites with minor allele frequency above 5%, suggesting that many of the common
286 variable sites have been sampled.

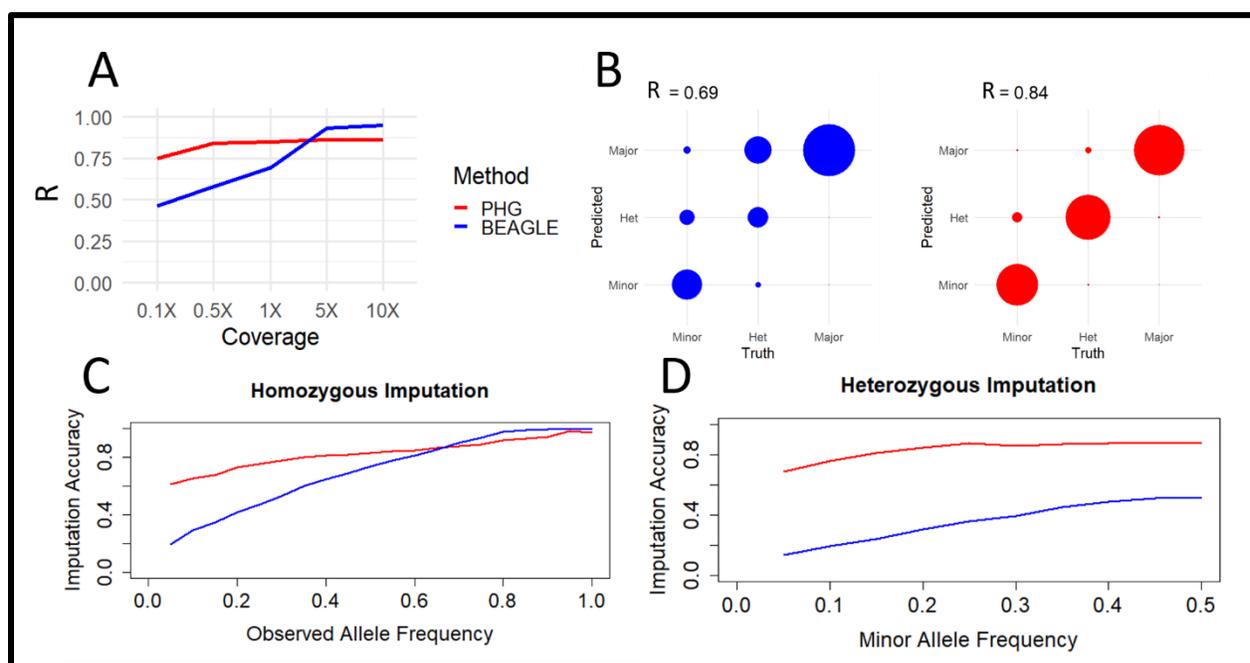
287 **Imputation and Genomic Prediction Accuracy**

288 Because imputation accuracy is dependent on the relative allele frequency and
289 phase of the allele being called, we classified genotype calls by allele frequency class:
290 homozygous major (both alleles are identical and have >50% allele frequency in
291 HapMapII), homozygous minor (both alleles are identical and have <50% allele
292 frequency in HapMapII), and heterozygous (two different alleles are present). In our
293 analyses, imputation accuracy is defined as the ability of the imputation method to
294 reconstitute genome-wide SNPs from the input data. We use the correlation between
295 the predicted alleles and the true alleles (defined by HapMapII) as a metric to make the
296 PHG and Beagle comparable, because the PHG utilizes reads and Beagle utilizes
297 variants to make their predictions.

298 Imputation of skim sequence genotyping showed PHG methods had a large
299 advantage over Beagle using low coverage sequence. At a level of 1X coverage
300 random sequencing, the PHG predicted allele calls with a correlation of $R^2=0.84$, while
301 the correlation between Beagle predicted alleles and the true calls was $R^2=0.69$ (Fig. 4
302 A). At higher depths of coverage (>5X), the raw data provides ample information to
303 distinguish between homozygous and heterozygous genotypes, allowing Beagle to
304 determine the correct genotype. The PHG, however, is able to distinguish between the
305 available haplotypes at a coverage of 0.5X and adding additional sequence data does

306 not increase the accuracy, as there is no correlation between accuracy and coverage
 307 beyond 0.5X.

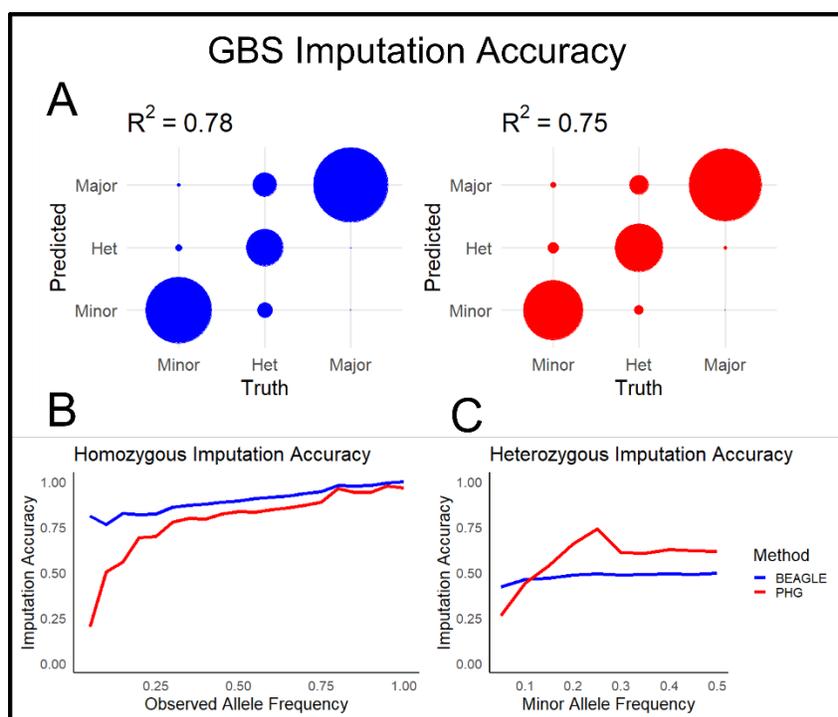
308 The improved performance of the PHG is most noticeable in its accurate
 309 predictions of heterozygous and rare genotypes. The PHG was able to impute
 310 genotypes with high accuracy regardless of allele class (Fig. 4B). The PHG's high
 311 accuracy at low allele frequencies for both homozygous (Fig. 4C) and heterozygous
 312 genotypes (Fig. 4D), display its ability to impute rare alleles.



313
 314 **Figure 4. Imputation Accuracy from skim sequencing. A) Displays correlation**
 315 **between imputed and true variants by imputing with the PHG and Beagle at**
 316 **different levels of skim sequencing. B) Displays concordance between true and**
 317 **imputed alleles at 1X coverage separated by alleles classes: minor, heterozygous,**
 318 **and major (circle radius is equal to the proportion of alleles in each class). C)**
 319 **Imputation accuracy at 1X coverage is shown for homozygous genotypes**
 320 **separated by allele frequency of the true allele at that locus. D) Imputation**

321 accuracy at 1X coverage is shown for heterozygous genotypes separated by
 322 minor allele frequency at that locus.

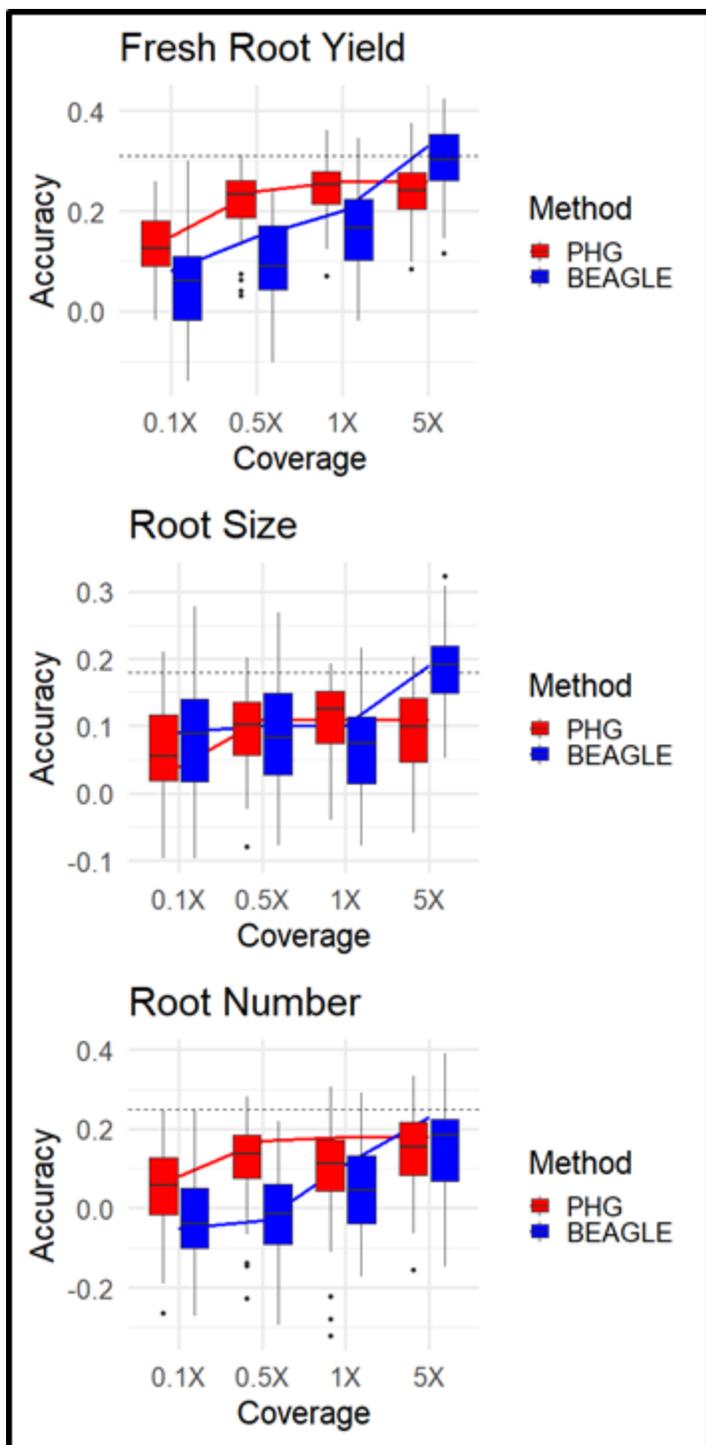
323 In addition to skim sequence scenarios, we also tested imputation using available
 324 GBS sequence for 20 clones. While skim sequence samples a random set of reads
 325 from across the genome, GBS is a replicable set of markers that a sparsely sampled
 326 across the genome. Imputation tests showed similar, but somewhat reduced
 327 accuracies using the PHG compared to Beagle (Fig. 5A). It is important to note
 328 however that the PHG still had improved accuracies in imputing heterozygous
 329 genotypes (Fig. 5C).



330
 331 **Figure 5. Imputation Accuracy from GBS sequencing. A) Displays concordance**
 332 **between true and imputed alleles separated by alleles classes (circle radius is**
 333 **equal to the proportion of alleles in each class) B) Imputation accuracy is shown**
 334 **for homozygous genotypes separated by allele frequency of the true allele at that**

335 **locus. C) Imputation accuracy is shown for heterozygous genotypes separated**
336 **by minor allele frequency at that locus.**

337 The imputed genotypes from skim sequence were then utilized in a genomic
338 prediction scheme consisting of 57 cassava clones (Supplemental Fig. 4) from a single
339 breeding program. Clones were selected from a single breeding program to minimize
340 confounding factors such as population structure and ensured an adequate level of
341 heritability to assess genomic prediction accuracy. Ten-fold cross validations and leave-
342 one-out validation showed that imputation accuracy generally appeared to follow the
343 trends in genomic prediction accuracy, for fresh root yield and root number, while no
344 clear pattern was apparent for the root size trait (Fig. 6).



345
 346 **Figure 6 Genomic Prediction Cross Validation. 10-Fold cross validation (box) and**
 347 **single holdout cross validation (line) show genomic prediction accuracies of 3**

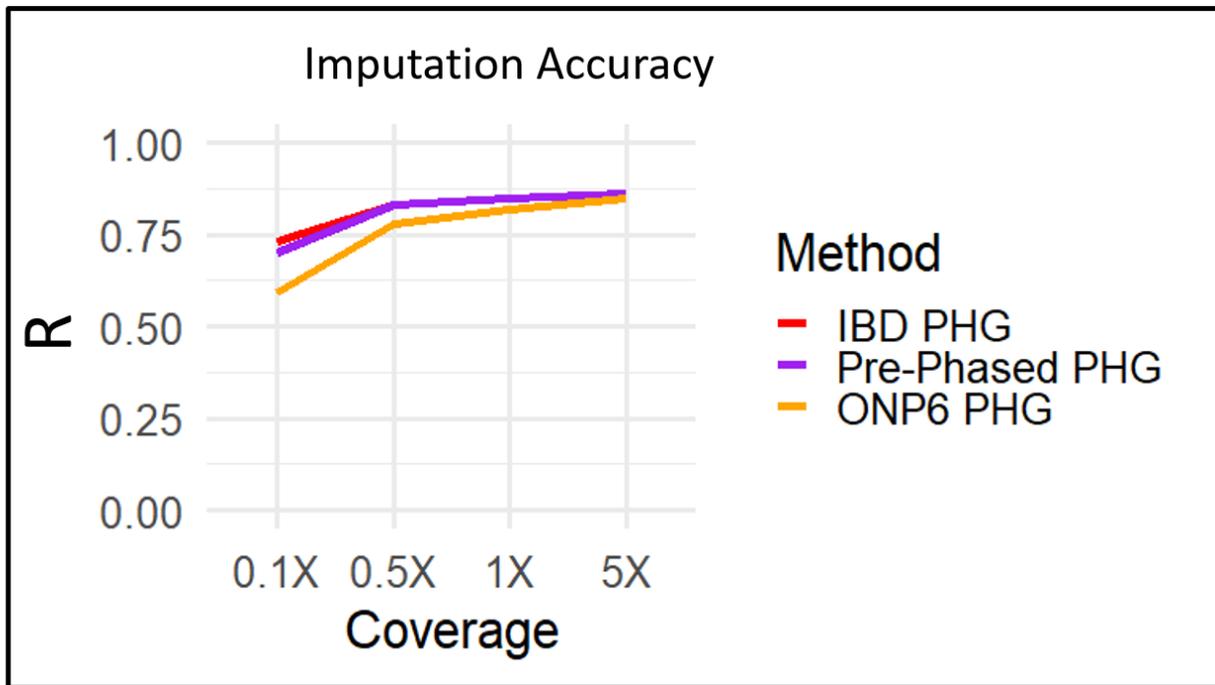
348 **root traits using different imputation methods at varied sequence depths. Single**
349 **holdout cross validation using complete genotype dataset is shown (dashed line).**

350

351 **Phased Haplotype PHG**

352 We tested the viability of populating the PHG with haplotypes phased by other
353 methods. We compared the IBD method of sampling phased haplotypes to two
354 methods of phasing variants. The first method used Beagle and HAPCUT2 to phase
355 the variants called from the HapMapII WGS data. The second method utilized 6
356 cassava clones with ONP long-read data. The IBD and Pre-Phased methods of
357 populating the cassava PHG produced almost identical accuracies (Fig. 7). These
358 results suggest that Beagle and HAPCUT could not accurately phase heterozygous
359 haplotypes at this scale, and the accurate haplotypes are derived from IBD haplotypes.
360 While the PHG was made from 6 clones with ONP data did perform as well as the other
361 methods, it relied on a far narrower set of germplasm. This suggests that accurate
362 haplotypes were likely captured using this method but lacked adequate sampling to
363 capture sufficient haplotypes.

364



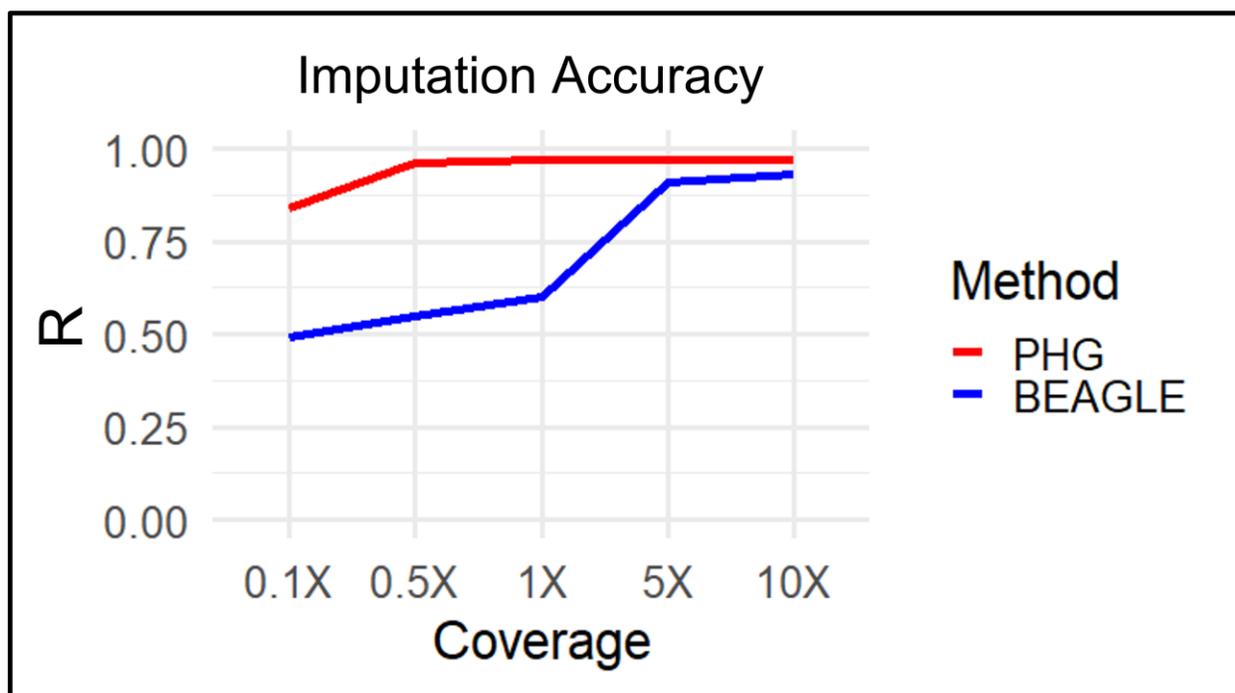
365
 366 **Figure 7. Haplotype Phasing Methods in the PHG. Imputation accuracy is shown**
 367 **for 3 different methods of populating a PHG. First the IBD PHG (red) was**
 368 **populated using homozygous haplotypes from the 241 HapMapII clones. Second,**
 369 **the Pre-Phased PHG (Purple) used Beagle and HPACUT2 to phase these same**
 370 **clones. Third, the ONP6 PHG (Yellow) used ONP long-reads and WhatsHap to**
 371 **phase 6 related taxa to the test set.**

372
 373 **Imputation Simulation**

374 Evident from the tests using haplotypes from IBD regions of the genome,
 375 sampling phased haplotypes is a difficult aspect of creating an effective PHG in a
 376 heterozygous species. To explore the performance of the PHG in a scenario where one
 377 could aptly sample the diversity of haplotypes, we used simulated offspring from a set of
 378 20 phased genomes. While phasing errors exist, we accepted these phases as truth for
 379 the simulation of offspring. This ensured that all haplotypes present in the offspring

380 exist in the PHG database. We found that the disparity in accuracies between PHG and
 381 Beagle at high sequence coverage disappeared in our simulation (Fig. 8), while the
 382 trend in Beagle accuracy was very similar to our empirical tests. While the simulation
 383 does represent an ideal scenario, including a narrower set of germplasm, it highlights
 384 the performance of the PHG when accurately phased haplotypes are available.

385



386

387 **Figure 8 Imputation Accuracy with Simulated Genotypes.** A simulated scenario
 388 where 20 parents with full phased information are used to populate a PHG.
 389 **Correlation between imputed and true variants by imputing with the PHG and**
 390 **Beagle at different levels of skim sequencing.**

391

392

393

394

DISCUSSION

395 We have detailed a method of implementing a PHG for the heterozygous plant
396 species cassava. This PHG database utilizes phased haplotypes to predict missing
397 genotypes from low depth input sequence. Runs of homozygosity formed by IBD
398 relationships proved to be a reliable method of sampling phased haplotypes given the
399 available data (Fig. 7). This method of obtaining haplotypes, while not able obtain the
400 full diversity of alleles, captured 77% of common alleles and produced ample
401 haplotypes for significant imputation accuracy at very low sequence depth (Fig. 4A).

402 The high accuracy of the PHG demonstrates its potential as an imputation tool
403 for use in heterozygous crops. The advantages of the PHG imputation methodology are
404 especially evident in its accuracy at calling rare and heterozygous alleles (Fig 4C,4D).
405 Furthermore, the observed weaker relationship between allele frequency and imputation
406 accuracy, highlights its ability to predict rare alleles. Across both simulated and
407 empirical experiments, we found that the ability of the PHG to impute whole genome
408 variants was consistent at or above 0.5X sequence coverage. The haplotype-based
409 representation of the genome enables this imputation methodology to overcome the
410 logistical hurdles such as those produced by sequencing and assembly errors, repetitive
411 sequences, and poor alignments.

412 The plateau reached in imputation accuracy (Fig. 4A) using the PHG most likely
413 indicates that we have not sufficiently sampled the diversity of possible haplotypes. At
414 sequence coverages of 5X and higher, the raw data can produce the true genotypes
415 and little imputation of missing genotypes is occurring. The PHG imputation is limited to
416 predicting haplotypes that are already present in the database, while Beagle can rely on
417 the genotypes called from the high depth (>1X) raw sequence, meaning that there is

418 much fewer missing data for Beagle to impute. This scenario of high depth sequence is
419 useful to diagnose challenges in imputation, however it does not correlate to many real
420 applications. The disparity between the PHG and Beagle at these high coverages
421 points to the presence of missing haplotypes in the database, rather than any disparity
422 in performance.

423 The hypothesis of missing haplotypes limiting imputation accuracy is supported
424 by a visible relationship between homozygous incidence in our population and reference
425 range imputation accuracy (Supplemental Fig. 5), suggesting that those ranges with
426 poor imputation accuracy were not amply sampled. The length and abundance of the
427 IBD runs of homozygosity in our dataset likely determine the ability of the HMM to
428 accurately predict haplotypes. There may be many factors that affect the prevalence of
429 IBD haplotypes including recessive deleterious effects, populations size, population
430 diversity, and heterozygosity. We saw that the disparity in imputation accuracy was
431 remedied under simulation, where all possible haplotypes were sampled in the
432 database (Fig. 8). These results suggest that, although an already powerful tool, the
433 PHG achieves maximum performance with sufficient sampling of available haplotypes.

434 Currently the performance using GBS data appears to be similar between the
435 PHG and Beagle (Fig. 5). Imputation from reduced representation genotyping such as
436 GBS is challenging due to the sparse sampling across the genome and varied levels of
437 sequence quality. Excellent imputation accuracy in inbred crops Sorghum (Jensen *et*
438 *al.* 2020) and Maize (Valdes Franco 2020) using these genotyping methods highlights
439 the potential benefits of the PHG in these scenarios. Because reduced representation
440 genotyping methods are likely the most commonly implemented, current efforts are

441 being made to improve heterozygous imputation using these technologies. We expect
442 improved haplotype sampling and phasing to improve imputation accuracy for these
443 genotyping platforms. Further haplotype sampling paired with developments in the
444 PHG imputation methodology will likely improve imputation accuracy from these
445 genotyping methods.

446 While the imputation accuracy of the PHG is limited based on the haplotype
447 sampling, its high accuracy with low levels of input sequence highlights its potential for
448 genomic applications, where sparse genotyping is common. We showed that this is
449 true regarding genomic prediction by performing cross-validations with the imputed
450 genotypes (Fig. 3). The genomic prediction was still limited by imputation accuracy, but
451 by enabling higher accuracy we can achieve more reliable predictions (Pimentel *et al.*
452 2015; Wang *et al.* 2016; Van Den Berg *et al.* 2017).

453 With increased imputation accuracy from more limited genotyping inputs, a
454 breeding program may be able to afford to cross and genotype more offspring, enabling
455 them to increase selection pressure across their breeding pool. Similarly, imputation to
456 genome-wide scale can bridge gaps between different data sets containing information
457 on different marker panels, enabling the use of larger datasets for prediction. Accurate
458 imputation could also enable breeders to utilize genomic prediction models that
459 incorporate more prior functional information on genome-wide variant effects into
460 predictions, using methods such as GFBLUP (Fang *et al.* 2017) or BayesR (MacLeod *et*
461 *al.* 2016; Van Den Berg *et al.* 2017). These possible applications of imputation have the
462 potential to increase total genetic gain made by breeding programs.

463 We show that while computational methods might not be able to solve haplotype
464 phasing with short-read data, long-read sequencing may be able to overcome that
465 issue. The Pre-Phased PHG produced similar accuracies to the IBD method,
466 suggesting the additional haplotypes added by phasing why heterozygous alleles using
467 Beagle and HAPCUT were not accurate over long distances. While limited in scope,
468 the ability of the PHG created from 6 clones with ONP data suggests the potential
469 application of long reads for obtaining phased haplotypes. One could envision a
470 breeding scenario in which parents are sequenced and phased using long-reads and
471 offspring are predicted from minimal genotyping input using the PHG. Then every few
472 generations shallow WGS can be used to update the PHG and compensate for
473 changing LD structures.

474 Applying the PHG to cassava and other heterozygous crops will depend on the
475 ability to sample phased haplotypes within the given population. We've shown that
476 utilizing high depth WGS data and IBD regions of the genome can be used to reliably
477 sample many phased haplotypes, and that the resulting PHG can impute with high
478 accuracy from low depth sequence. This method of sampling haplotypes will be highly
479 dependent on the diversity and heterozygosity of the species and population for any
480 given application. Other necessary considerations for the decision to use the PHG
481 include genome size, reference genome quality, training data availability, species
482 ploidy. In applications where imputation is commonly implemented, training data that
483 can be used to construct a PHG may already be available. Our long-read data results
484 show the potential for more easily capturing phased haplotypes as these technologies
485 become more available. Using genome assemblies produced from long-reads as inputs

486 to the PHG has been shown to be very effective in Maize, while this method has not
487 been implemented in outbred species. The potential for the PHG as a tool in
488 heterozygous crops has been shown here, while the specific methods to produce the
489 phased haplotypes will have to be designed around the target species and scenario.

490 CONCLUSION

491 The PHG is a method to reduce a genome to a set of haplotypes. We have
492 shown that this method can be used to predict whole genome haplotypes in a
493 heterozygous species from sparse genotyping information. Its high accuracy, especially
494 in rare alleles, at very low depths of skim sequencing makes it a potentially powerful
495 imputation tool. Continued work in populating the PHG database with confidently
496 phased haplotypes will lead to a more consistent prediction model across varied
497 genotyping methods.

498 DATA AVAILABILITY

499 Supplementary files and scripts used for the production and testing of the cassava PHG
500 can be found at https://bitbucket.org/bucklerlab/p_cassava_phg. Genotype and
501 phenotype data from HapMapII (Ramu *et al.* 2017) was downloaded from
502 cassavabase.org. Support and methods for practical haplotype graph implementation
503 can also be found at <https://bitbucket.org/bucklerlab/practicalhaplotypegraph/wiki/Home>.
504 Raw Oxford nanopore sequence data for this project is available at NCBI BioProject ID
505 PRJNA589272.

506 ACKNOWLEDGMENTS

507 We'd like to acknowledge the programming staff in the Buckler lab who created
508 and support the development of the Practical Haplotype Graph, as well as other lab

509 members that provided feedback on experimental design. We are also grateful for the
510 greater Nextgen cassava community for supporting the curation of genotype and
511 phenotype data used in this project as well as the organization of this data in
512 cassavabase.org.

513 FUNDING

514 This study is made possible by the funding and support of the Nextgen Cassava project,
515 the Bill and Malinda Gates foundation, and the USDA-ARS.

516 COMPETING INTERESTS

517 The authors declare no competing financial interests.

518

519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541

REFERENCES

Alipour, H., G. Bai, G. Zhang, M. R. Bihamta, V. Mohammadi *et al.*, 2019 Imputation accuracy of wheat genotyping-by-sequencing (GBS) data using barley and wheat genome references. *PLoS One* 14:.

Van Den Berg, I., P. J. Bowman, I. M. MacLeod, B. J. Hayes, T. Wang *et al.*, 2017 Multi-breed genomic prediction using Bayes R with sequence data and dropping variants with a small effect. *Genet. Sel. Evol.* 49: 1–15.

Browning, B. L., Y. Zhou, and S. R. Browning, 2018 A One-Penny Imputed Genome from Next-Generation Reference Panels. *Am. J. Hum. Genet.* 103: 338–348.

Cleveland, M. A., J. M. Hickey, and B. P. Kinghorn, 2011 Genotype imputation for the prediction of genomic breeding values in non-genotyped and low-density genotyped individuals, pp. S6 in *BMC Proceedings*, BioMed Central.

Edge, P., V. Bafna, and V. Bansal, 2017 HapCUT2: Robust and accurate haplotype assembly for diverse sequencing technologies. *Genome Res.* 27: 801–812.

Fang, L., G. Sahana, P. Ma, G. Su, Y. Yu *et al.*, 2017 Exploring the genetic architecture and improving genomic prediction accuracy for mastitis and milk production traits in dairy cattle by mapping variants to hepatic transcriptomic regions responsive to intra-mammary infection. *Genet Sel Evol* 49: 44.

Fragoso, C. A., C. Heffelfinger, H. Zhao, and S. L. Dellaporta, 2016 Imputing Genotypes in Biallelic Populations from Low-Coverage Sequence Data. *Genetics* 202: 487–495.

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

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

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

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

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

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

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

556 Kono, T. J. Y., C. Liu, E. E. Vonderharr, D. Koenig, J. C. Fay *et al.*, 2019 The Fate of
557 Deleterious Variants in a Barley Genomic Prediction Population. Genetics 213:
558 1531–1544.

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

562 Loh, P. R., P. F. Palamara, and A. L. Price, 2016 Fast and accurate long-range phasing
563 in a UK Biobank cohort. Nat. Genet. 48: 811–816.

564 MacLeod, I. M., P. J. Bowman, C. J. Vander Jagt, M. Haile-Mariam, K. E. Kemper *et al.*,

565 2016 Exploiting biological priors and sequence variants enhances QTL discovery
566 and genomic prediction of complex traits. *BMC Genomics* 17: 1–21.

567 Meuwissen, T. H. E., B. J. Hayes, and M. E. Goddard, 2001 Prediction of total genetic
568 value using genome-wide dense marker maps. *Genetics* 157: 1819–1829.

569 Nazzicari, N., F. Biscarini, P. Cozzi, E. C. Brummer, and P. Annicchiarico, 2016 Marker
570 imputation efficiency for genotyping-by-sequencing data in rice (*Oryza sativa*) and
571 alfalfa (*Medicago sativa*). *Mol. Breed.* 36: 69.

572 Pimentel, E. C. G., C. Edel, R. Emmerling, and K. U. Götz, 2015 How imputation errors
573 bias genomic predictions. *J. Dairy Sci.* 98: 4131–4138.

574 Ramu, P., W. Esuma, R. Kawuki, I. Y. Rabbi, C. Egesi *et al.*, 2017 Cassava haplotype
575 map highlights fixation of deleterious mutations during clonal propagation. *Nat.*
576 *Genet.* 49: 959–963.

577 Romay, M. C., 2018 Rapid, Affordable, and Scalable Genotyping for Germplasm
578 Exploration in Maize, pp. 31–46 in Springer, Cham.

579 Schrunner, S. D., R. S. Mari, J. Ebler, M. Rautiainen, L. Seillier *et al.*, 2020 Haplotype
580 Threading: Accurate Polyploid Phasing from Long Reads. *bioRxiv*
581 2020.02.04.933523.

582 Spencer, C. C. A., Z. Su, P. Donnelly, and J. Marchini, 2009 Designing Genome-Wide
583 Association Studies: Sample Size, Power, Imputation, and the Choice of
584 Genotyping Chip (J. D. Storey, Ed.). *PLoS Genet.* 5: e1000477.

585 Stephens, Z. D., M. E. Hudson, L. S. Mainzer, M. Taschuk, M. R. Weber *et al.*, 2016
586 Simulating next-generation sequencing datasets from empirical mutation and
587 sequencing models. *PLoS One* 11:.

588 Swarts, K., H. Li, J. A. Romero Navarro, D. An, M. C. Romay *et al.*, 2015 Novel
589 Methods to Optimize Genotypic Imputation for Low-Coverage, Next-Generation
590 Sequence Data in Crop Plants. *Plant Genome* 7: 0.

591 Torkamaneh, D., B. Boyle, and F. Belzile, 2018 Efficient genome-wide genotyping
592 strategies and data integration in crop plants. *Theor. Appl. Genet.* 131: 499–511.

593 Valdes Franco, J. A., 2020 A Maize Practical Haplotype Graph Leverages Diverse NAM
594 Assemblies. *Zachary R. Mill.* 2: 0.

595 Wang, Y., G. Lin, C. Li, and P. Stothard, 2016 Genotype Imputation Methods and Their
596 Effects on Genomic Predictions in Cattle. *Springer Sci. Rev.* 4: 79–98.

597 Xu, Y., X. Liu, J. Fu, H. Wang, J. Wang *et al.*, 2020 Enhancing Genetic Gain through
598 Genomic Selection: From Livestock to Plants. *Plant Commun.* 1: 100005.

599 Yang, J., S. Mezmouk, A. Baumgarten, E. S. Buckler, K. E. Guill *et al.*, 2017 Incomplete
600 dominance of deleterious alleles contributes substantially to trait variation and
601 heterosis in maize (J. C. Fay, Ed.). *PLOS Genet.* 13: e1007019.

602 Yun, L., C. Willer, S. Sanna, and G. Abecasis, 2009 Genotype imputation. *Annu. Rev.*
603 *Genomics Hum. Genet.* 10: 387–406.

604

605