

# Phased genotyping-by-sequencing enhances analysis of genetic diversity and reveals divergent copy number variants in maize

Heather Manching<sup>\*</sup>, Subhajit Sengupta<sup>†</sup>, Keith R. Hopper<sup>‡</sup>, Shawn W. Polson<sup>§</sup>, Yuan Ji<sup>†,\*\*</sup> and Randall J. Wisser<sup>\*,1</sup>

<sup>\*</sup>Department of Plant and Soil Sciences, University of Delaware, Newark, DE 19716, <sup>†</sup>Program of Computational Genomics & Medicine, NorthShore University HealthSystem, Evanston, IL 60201, <sup>‡</sup>Beneficial Insect Introductions Research Unit, United States Department of Agriculture, Agricultural Research Service, Newark, DE 19713, <sup>§</sup>Center for Bioinformatics and Computational Biology, Delaware Biotechnology Institute, University of Delaware, Newark, DE 19711, <sup>\*\*</sup>Department of Health Studies, University of Chicago, Chicago, IL 60637

**ABSTRACT** High-throughput sequencing of reduced representation genomic libraries has ushered in an era of genotyping-by-sequencing (GBS), where genome-wide genotype data can be obtained for nearly any species. However, there remains a need for imputation-free GBS methods for genotyping large samples taken from heterogeneous populations of heterozygous individuals. This requires a number of issues encountered with GBS be considered, including the sequencing of non-overlapping sets of loci across multiple GBS libraries, a common missing data problem that results in low call rates for markers per individual, and a tendency for applicability only in inbred line samples with sufficient linkage disequilibrium for accurate imputation. We addressed these issues while developing and validating a new, comprehensive platform for GBS. This study supports the notion that GBS can be tailored to particular aims, and using *Zea mays* our results indicate that large samples of unknown pedigree can be genotyped to obtain complete and accurate GBS data. Optimizing size selection to sequence a high proportion of shared loci among individuals in different libraries and using simple *in silico* filters, a GBS procedure was established that produces high call rates per marker (>85%) with accuracy exceeding 99.4%. Furthermore, by capitalizing on the sequence-read structure of GBS data (stacks of reads), a new tool for resolving local haplotypes and scoring phased genotypes was developed, a feature that is not available in many GBS pipelines. Using local haplotypes reduces the marker dimensionality of the genotype matrix while increasing the informativeness of the data. Phased GBS in maize also revealed the existence of reproducibly inaccurate (apparent accuracy) genotypes that were due to divergent copy number variants unobservable in the underlying single nucleotide polymorphism data.

## KEYWORDS

GBS; haplotype phasing; copy number variant; imputation; maize

1

## 2 INTRODUCTION

3 Genome-wide genotyping of population samples is fundamental  
4 to a range of studies in genetics and genomics, and genotyping-  
5 by-sequencing (GBS) of multiplexed, high-throughput sequencing  
6 (HTS) libraries has emerged as a cost-effective strategy for ob-  
7 taining this information. Genotyping-by-sequencing of multiple  
8 samples relies on a reduced representation sequencing strategy

9 that restricts sequencing across what is hoped to be a common sub-  
10 space of the genome in different samples. A prevalent technique  
11 used for this involves the ligation of barcoded adapters to DNA  
12 digested with restriction enzymes, followed by sequencing of frag-  
13 ments within a restricted size range. In principle, this leads to  
14 stacks of sequences anchored at the restriction cut sites across the  
15 genomes of different individuals. Following initial publications on  
16 restriction enzyme-mediated GBS (e.g. Baird *et al.* 2008; Andolfatto  
17 *et al.* 2011; Elshire *et al.* 2011), the approach has continued to be  
18 extended and optimized, such as fine-tuning the number of loci  
19 sequenced by using different restriction enzymes and size selection

Copyright © 2017 Manching *et al.*

Manuscript compiled: May 2017

<sup>1</sup>Randall J. Wisser; Department of Plant and Soil Sciences, University of Delaware, 152 Townsend Hall, 531 South College Avenue, Newark, DE, 19716; rjw@udel.edu

20 windows (Peterson *et al.* 2012) or selective primers (Sonah *et al.* 21 2013), minimizing the front-end cost by using enzymes that create 22 blunt-end fragments for universal adapters (Heffelfinger *et al.* 23 2014), and increasing the scalability for large samples by including 24 a sequence capture step (Ali *et al.* 2016). Methods and applications 25 of GBS continue to develop at a rapid rate.

26 With the mass of data produced by HTS, even low sequencing 27 error rates can lead to an abundance of inaccurate genotype 28 calls. Moreover, errors accumulate along steps in the GBS pipeline 29 from library preparation to data processing. As a reflection of the 30 amount of noise in GBS data, it is not unusual to find publications 31 (as well as our own experience) where as much as ~50% of the 32 raw HTS data is filtered after mapping (e.g. Beissinger *et al.* 2013) 33 and as much as ~95% of the discovered variants are filtered after 34 variant calling (e.g. Hyma *et al.* 2015). Typically, statistically 35 or biologically informed criteria for filtering are used to enrich 36 for more accurate GBS data, or a model is fitted to the data to 37 remove genotypes that do not conform to expectations (e.g. tests for 38 independent assortment and co-segregation in genetic mapping 39 populations, Hardy-Weinberg equilibrium in other population 40 samples). However, estimates of genotyping accuracy for a given 41 protocol, pipeline and application of GBS (under all of its flavors: 42 e.g. RAD, GBS, ddRAD, Rapture, etc.) are rarely reported. One 43 recent study compared different bioinformatic pipelines and found 44 that genotyping accuracy for one GBS dataset ranged from 76% 45 to 99%, depending on the bioinformatic pipeline used to score 46 genotypes (Torkamaneh *et al.* 2016). Accuracy estimates based 47 on resequencing of GBS loci have ranged from 51% (Rocher *et al.* 48 2015) to 99% (Torkamaneh *et al.* 2016). Given the wide range in 49 accuracy and method- and study-specific nature of GBS, establish- 50 ing controls or approaches that allow accuracy to be estimated is 51 important when embarking on GBS studies.

52 Genotyping-by-sequencing tends to have a missing-data prob- 53 lem due to two primary issues: (i) for libraries with low sequence 54 coverage per sample, the amount of missing genotype data can be 55 high, resulting in a low call rate per marker (as low as 10%; Fu and 56 Peterson 2011); (ii) for large population samples where many li- 57 braries need to be sequenced, the number of shared loci in different 58 libraries may be low, resulting in low call rates per sample. This lat- 59 ter issue has not been fully addressed in the literature on GBS. We 60 are aware of one study in which the consistency of genotyping the 61 same loci in different libraries with the same samples was assessed, 62 and in this study, overlap of loci among libraries ranged from 12% 63 to 98% (DaCosta and Sorenson 2014). The results of that study 64 suggested an impractical (costly) solution for genotyping large 65 samples, where greater consistency would be achieved by pooling 66 samples later in the protocol. This would require individual sam- 67 ples to be processed through size selection and PCR amplification 68 independently, which increases the cost for reagents and time for 69 library preparation. A solution to the GBS missing-data problem 70 has been imputation, whereby the genotype for a missing SNP is 71 inferred from the state of nearby SNPs. However, imputation is not 72 always possible and may not be sufficiently accurate, in which case 73 both of the above issues compromise the effectiveness of GBS in 74 studies where sample sizes exceed the multiplex depth of a single 75 sequencing library. While the former issue may be addressed by 76 reducing the sequence space (e.g. Peterson *et al.* 2012; Sonah *et al.* 77 2013) or deeper sequencing, the latter issue requires techniques 78 that provide reproducible enrichment of shared loci (e.g. Ali *et al.* 79 2016).

80 Despite the shortcomings mentioned above, GBS has been ef- 81 fective in several settings. There is now a growing interest in using

82 GBS for a wider range of studies (Andrews *et al.* 2016), many of 83 which could benefit from or be advanced with phased genotype 84 data. The standard approach when performing GBS involves scor- 85 ing single nucleotide polymorphisms (SNPs). However, in GBS 86 data, locally phased haplotypes exist as stacks of reads for each 87 locus. If these contain more than one SNP, then multi-nucleotide 88 polymorphisms (MNPs) can be genotyped. The use of MNPs has 89 not been widely exploited in GBS due to the lack of tools for scoring 90 phased genotypes. The software STACKS (Catchen *et al.* 2013) and 91 Haplotag (Tinker *et al.* 2016) do have functions embedded within 92 their pipelines for extracting MNP haplotypes from GBS data, but 93 the pipeline-dependency of these functions limits their use more 94 generally, and the algorithms rely on population-specific filtering 95 criteria rather than statistical evaluation of the likelihood that each 96 haplotype in an individual is real. Because HTS data is imperfect 97 and the quality of sequenced bases and read mappings are quanti- 98 tatively encoded, scoring MNPs at a locus within an individual is 99 not straightforward if one considers this information relevant. A 100 local haplotyping tool, LocHap, which uses community-standard 101 file formats, offers a more generalized solution for phased geno- 102 typing. LocHap was developed under a probabilistic framework 103 that uses the quality metrics of base calls and mapping results as 104 well as sampling effects to phase MNPs in HTS data (Sengupta 105 *et al.* 2015). LocHap was designed to identify distinct haplotypes in 106 heterogeneous populations of cells (somatic mosaicism) irrespec- 107 tive of homology, such that the outputted haplotypes at a single 108 locus can vary in length and are not necessarily alignable within 109 and across individuals. Consequently, LocHap does not produce 110 data in a format that is compatible with population or quantitative 111 genetic studies. In our study, LocHap was extended to perform 112 phased genotyping with GBS data on individual samples, which 113 we call LocHap-GBS.

114 Phased genotyping is useful for various applications in genet- 115 ics and genomics. Typing MNPs can help distinguish more than 116 two alleles at a locus, providing greater information content for 117 studying genetic diversity (Lu *et al.* 2011). Phased genotypes can 118 facilitate imputation in multi-parental populations where SNPs 119 would otherwise conflate ancestral alleles (Davies *et al.* 2016). Hap- 120 lotype data can also increase the accuracy of estimated breeding 121 values, thereby increasing the efficiency of plant and animal breed- 122 ing (Ferdosi *et al.* 2016). In our study, while validating LocHap-GBS, 123 we found that phased genotyping can also uncover copy number 124 polymorphisms.

125 The aims of this study were to: (i) establish a standardized and 126 empirically optimized flex-plex GBS protocol with an accompany- 127 ing informatics pipeline for genotyping; (ii) evaluate the accuracy 128 and potential applicability of this procedure for genotyping large 129 samples from heterogeneous populations of heterozygous individ- 130 uals; and (iii) extend the use of GBS for phased genotyping. We 131 used maize, which is an agriculturally relevant species with an 132 ~2.3 Gb completed reference genome sequence, but where genomic 133 analysis is challenged by large amounts of repetitive sequence. 134 Genetic trios of pairs of inbred parental lines and their F<sub>1</sub> progeny 135 were used in assessing accuracy of GBS. We tested whether simple 136 *in silico* filters could facilitate highly accurate scoring of genotypes 137 irrespective of zygosity and knowledge of parentage. In addition 138 to developing and validating a method for phased-GBS, we also 139 examined how SNP and MNP data affect inference on the relat- 140 edness among a set of inbred lines used by the maize genetics 141 research community.

143 **Study Samples**

144 Detailed information on the samples used in this study is in Supple-  
 145 mentary Material (File S1). Briefly, the samples included a synthetic  
 146 population created from seven tropical inbred lines (used to evalu-  
 147 ate the consistency of sequenced loci in different GBS libraries),  
 148 sets of trios or pairs of inbred lines and their corresponding F<sub>1</sub>  
 149 hybrids (used to evaluate the repeatability and accuracy of GBS),  
 150 an F<sub>2</sub> population (used to examine the transmission of phased  
 151 MNPs, along with the corresponding parental trio), and parental  
 152 inbred lines of a maize nested association mapping (NAM) popu-  
 153 lation (McMullen *et al.* 2009; used to examine phased genotyping  
 154 for analysis of genetic diversity).

155 **Genotyping-by-Sequencing**

156 We describe the design, protocol, and associated software for GBS  
 157 based on a double digestion technique similar to protocols by  
 158 Poland *et al.* (2012) and Peterson *et al.* (2012). A specific interest of  
 159 ours was to develop a GBS design for genotyping heterozygous  
 160 samples of potentially unknown parentage with little missing data  
 161 for large sample studies. The method was optimized, validated,  
 162 and tested for various applications.

163 Detailed protocols and other relevant information are pro-  
 164 vided in Supplemental Material (File S2, S3). Briefly, two sep-  
 165 arate but related restriction-associated sequence polymorphism  
 166 (RASP) adapter designs were developed: (i) RASP-1.0 was based  
 167 on Illumina's "genomic DNA" adapters (Illumina, Inc., San Diego,  
 168 CA) modified with appropriate overhang sequences for ligation  
 169 and a 6-nucleotide inline barcode for 48-plex genotyping (RASP-  
 170 1.0 and RASP-1.1 adapter oligonucleotides and primers; File S3);  
 171 (ii) RASP-2.0 was based on Illumina's TruSeq adapters modi-  
 172 fied with appropriate sequence overhangs for ligation. RASP-  
 173 2.0 uses variable length inline barcodes (5-10 nucleotides) for 48-  
 174 plex genotyping (barcode sequences were designed using <http://www.deenabio.com/gbs-adapters>) along with standard TruSeq in-  
 175 dices. The inline barcoded adapters and TruSeq barcoded adapters  
 176 can be combined to construct plexes in multiples of 48 (RASP-2.0  
 177 adapter oligos and primers: File S3).

179 The protocol used for GBS was improved over the course of  
 180 our study, including changes to the adapter sets, the size selection  
 181 method and number of reactions used for PCR (Table S1). The gen-  
 182 eral protocol was as follows. Purified and normalized DNA (200  
 183 ng) was digested using two restriction endonucleases, *NgoMIV*  
 184 and *Csp6I*, for 30 minutes at 37 °C. The *NgoMIV/Csp6I* enzyme  
 185 pair was chosen because it provided the highest read output and  
 186 the lowest variation between samples based on comparisons be-  
 187 tween four different enzyme pairs (*NgoMIV/Csp6I*; *NgoMIV/MseI*;  
 188 *PstI/Csp6I*; *PstI/MseI*; data not shown). Adapters were ligated  
 189 with T4 DNA ligase using temperature-cycle ligation (Lund *et al.*  
 190 1996; 300 cycles of 30s at 10 °C and 30s at 30 °C; this was deter-  
 191 mined to improve ligation efficiency by qPCR [data not shown]),  
 192 followed by heat inactivation of the ligase at 65 °C for 30 minutes.  
 193 An equal volume of each ligate was pooled and then purified using  
 194 either AMPure (A63880; Beckman Coulter, Inc., CA, USA) or SPRI-  
 195 elect (B23317; Beckman Coulter, Inc., CA, USA) beads following  
 196 the manufacturer's recommended protocol for standard clean-up  
 197 (exclusion of fragments <100 bp). For all libraries, prior to PCR am-  
 198 plification, size selection was performed on the pooled ligate using  
 199 a BluePippin (Sage Science, MA, USA). Size selected samples were  
 200 PCR amplified with Phusion High Fidelity Master Mix (M0531S;  
 201 NEB, Inc., MA, USA) using the universal primer sequences from  
 202 Illumina's genomic DNA or TruSeq sample kit. To reduce PCR bias,

203 a minimum of eight separate PCR reactions were performed on the  
 204 size selected template and the products were pooled. AMPure or  
 205 SPRI beads were used to eliminate excess primers and bi-products  
 206 (exclusion of fragments <100 bp). The range and peak of Bioana-  
 207 lyzer size fragment profiles were analyzed for each library before  
 208 and after PCR to determine the consistency of the size profiles. Due  
 209 to variation in fragment length profiles found following pre-PCR  
 210 size selection, a post-PCR size selection was performed on some of  
 211 the libraries, followed by Bioanalyzer analyses to confirm size pro-  
 212 files. This turned out to be critical for genotyping separate libraries,  
 213 and we have since included this change in our standard protocol.  
 214 Multiplex libraries were quantified using Quant-iT PicoGreen ds-  
 215 DNA assay kit (P7589; Thermo Fisher Scientific, Inc., MA, USA)  
 216 and sequenced (1 x 101 cycles) on an Illumina HiSeq 2500 at the  
 217 Delaware Biotechnology Institute .

218 **GBS data processing**

219 **Computational pipeline:** Sequences were processed using a cus-  
 220 tom *reduced representation* "RedRep" computational pipeline  
 221 with the following basic steps: 1) sequence splitting by bar-  
 222 code; 2) raw sequence quality control; 3) reference genome map-  
 223 ping; and 4) variant calling (SNPs). Briefly, sequences are decon-  
 224 voluted by barcode using custom logic and the FASTX-Toolkit  
 225 ([http://hannonlab.cshl.edu/fastx\\_toolkit](http://hannonlab.cshl.edu/fastx_toolkit)). Novel scripts and the Cu-  
 226 tAdapt package (Martin 2011) are then used to remove adapters,  
 227 trim low quality read ends, and filter out sequences that do not  
 228 meet minimum length/quality standards or lack the expected re-  
 229 striction site and adapter sequence expectations. Quality reads are  
 230 mapped to the reference genome using the BWA-MEM algorithm  
 231 (Li 2013) and SNPs are identified using the multi-sample discovery  
 232 mode of the GATK Unified Genotyper (McKenna *et al.* 2010). The  
 233 scripts for the pipeline and documentation are available under the  
 234 open source MIT license at <https://github.com/UD-CBCB/RedRep>.  
 235 Version 2.0 scripts committed to the repository on December 7,  
 236 2016 were used for this manuscript.

237 **Post-vcf filtering:** All 64-384 bp loci flanked by *NgoMIV*  
 238 (G↓CCGGC) and *Csp6I* (G↓TAC) recognition sites in the B73 v2 re-  
 239 ference genome were identified by *in silico* digestion. BLASTn was  
 240 used to perform genome-wide searches for each *in silico* digested  
 241 sequence starting from the *Csp6I* recognition site (starting point  
 242 for sequence reads) up to a maximum of 96 nucleotides (longest  
 243 quality trimmed read length of our actual data). Sequences that  
 244 were more than four percent distant from all other loci in the refer-  
 245 ence genome were flagged as "4PD" loci, and SNPs that occurred  
 246 within these loci were maintained while the remaining SNPs were  
 247 filtered. This excludes the possibility for scoring variants at loci  
 248 within the reference genome that differ from one another by fewer  
 249 than four nucleotides in 100 base pairs. In addition, for each in-  
 250 dividual and at each SNP site, if the genotype was not based on  
 251 a minimum depth of read coverage, it was set to missing: 3X for  
 252 scoring inbred samples; 12X for scoring heterozygous samples.  
 253 Samples were then removed that had genotype data for <25% of  
 254 the loci for 48-plex libraries and <10% for 192-plex libraries; these  
 255 included problematic samples and negative controls. Finally, for  
 256 each SNP, an 85% call rate threshold was used, and SNPs with  
 257 more than two nucleotide variants were removed.

258 **In silico characterization of GBS loci:** The *in silico* digested loci,  
 259 along with the SNP sites scored across all samples were sum-  
 260 marized according their distribution across the B73 v2 reference  
 261 genome. Using the closest-features program of BEDOPS v2.4.15  
 262 (Neph *et al.* 2012), genic and intergenic associations of the digestion

263 loci and SNPs were determined using gene structures in release-5b  
264 (filtered gene set; <http://ftp.maizesequence.org/release-5b/>).

### 265 Assessing Accuracy of GBS

266 A set of parent-hybrid trios were used to assess scoring accuracy: (i)  
267 CML373, CML341, and CML373 X CML341; (ii) CML341, CML277,  
268 and CML341 X CML277; (iii) Tzi9, CML258, and Tzi9 X CML258.  
269 One trio involving Tzi8 was excluded from accuracy assessment be-  
270 cause of excess heterogeneity (~12.8%) and relatively high residual  
271 heterozygosity (~1.3%) as determined in preliminary work using  
272 the MaizeSNP50 BeadChip (Ganal *et al.* 2011). The repeatability  
273 of genotyping was estimated from replicate DNA samples of each  
274 member of the trio processed in the same library (except for the  
275 hybrid Tzi9 X CML258 because of quality issues with the repli-  
276 cate). Using loci with complete and consistent calls (i.e. between  
277 replicated DNA samples), genotyping accuracy was measured as  
278 the proportion of loci with the expected genotype in the hybrid  
279 given the genotype of the parents. Because GBS was performed  
280 on only one plant for each member of a trio, loci that were het-  
281 erozygous in either parental line were excluded when estimating  
282 genotyping accuracy (according to prior data of ours based on the  
283 MaizeSNP50 chip, residual heterozygosity and heterogeneity for the  
284 parental lines was <0.8%). Furthermore, because the one F<sub>1</sub>  
285 plant genotyped was taken from bulked seed of progeny from mul-  
286 tiple crosses between the parental lines (such that the genotyped  
287 parents may not provide the exact expectation for the specific F<sub>1</sub>  
288 plant that was genotyped), estimates of accuracy are expected to  
289 be slightly downward biased.

### 290 Genotyping of local haplotypes

291 Because GBS can resolve multiple variants present within each  
292 position-specific stack of sequence reads, phased haplotypes can  
293 be extracted. A local haplotyping program called LocHap (Sen-  
294 gupta *et al.* 2015) was extended for use with GBS data, which we  
295 refer to as LocHap-GBS (<http://compgenome.org/lochap/GBS/>). Fig-  
296 ure 1 depicts the workflow for LocHap-GBS, which requires four  
297 input files: the standard *bam*, *bai*, and *vcf* files, plus a user-specified  
298 *bed* file containing intervals where haplotypes should be searched.  
299 To develop LocHap-GBS, a collection of modules were created that  
300 provide the following functionalities: (i) a file parser that defines  
301 intervals in which to search for haplotypes (ii) automation for run-  
302 ning multiple samples and merging samples into a single output  
303 file; (iii) outputting of haplotypes that include bases even if they  
304 are homozygous in an individual; and (iv) additional flexibility in  
305 the format of output (either haplotype calling format [*hcf*] or *bed*).

306 While developing LocHap-GBS some additional improvements  
307 were made to LocHap that were related to dynamic memory allo-  
308 cation. Presently, LocHap is limited to calling haplotypes across a  
309 maximum of three heterozygous sites. LocHap-GBS uses a *bed* file  
310 with predefined intervals of windows in which to search for haplo-  
311 types, but these windows may contain more than three heterozy-  
312 gous sites in some individuals. Therefore, LocHap-GBS parses  
313 windows into subwindows that include a maximum of three het-  
314 erozygous SNPs by referencing the coordinates of each window  
315 and the genotypes of SNPs in the *vcf* file (Figure 1). In this study,  
316 the intervals specified in the *bed* file were the 4PD loci from the *in*  
317 *silico* digest (only SNPs at these loci had been maintained in our *vcf*  
318 outputs, so this was the genotyping space of interest in this study).

319 We define SNPs, MNPs, and copy number polymorphisms  
320 (CNPs) as the property of a locus (loci that have variation among in-  
321 dividuals) and single nucleotide variants (SNVs), multi-nucleotide  
322 variants (MNVs), and copy number variants (CNVs) as the prop-

erty of an individual or in reference to a specific form of the poly-  
324 morphism (a diploid individual is expected to have a maximum of  
325 two SNVs or MNVs at a locus). To prepare MNP data for analy-  
326 sis, MNVs were split into separate columns using the *df2genind*  
327 function within the *adegenet* package (Jombart 2008) in R (R Core  
328 Team 2016). An MNV containing one missing SNV was assigned  
329 a value of "NA", and a genotype containing a missing MNV was  
330 assigned a value of "NA."

### 331 Analysis of Genetic Diversity

332 To examine differences in information content when using SNP  
333 versus MNP data, genetic distances of the parents of the NAM  
334 population were analyzed. When analyzing the MNP data, we  
335 encountered genotypes with more than two MNVs (implying a  
336 ploidy of greater than two, which is not possible for orthologous  
337 loci in maize). After determining these were putative CNPs within  
338 the genomes of the sequenced samples relative to the reference  
339 genome sample (see Results), we examined the impact of includ-  
340 ing or excluding CNPs in the analysis of genetic diversity. The  
341 original MNP genotype matrix was split into three datasets where:  
342 (i) all putative CNPs were included for analysis (MNP1); (ii) all  
343 loci that had more than two MNVs were masked (MNP2); and (iii)  
344 loci that had more than two MNVs or were heterozygous in an  
345 inbred line (these are also potentially CNPs) were masked (MNP3).  
346 Pairwise genetic distance matrices based on shared allele distances  
347 (Bowcock *et al.* 1994) were computed for each dataset. The dis-  
348 tance matrices were compared based on summary statistics and the  
349 Mantel test for correlation between matrices (Mantel 1967; imple-  
350 mented using the *ade4* package in R). To visualize the relationships  
351 between the lines, multidimensional scaling (MDS, Kruskal 1964)  
352 was performed using the *cmdscale* function in the R package *stats*.  
353 Phylogenetic trees were generated using the BioNJ algorithm (Gas-  
354 cuel 1997) with the *ape* package (Paradis *et al.* 2004) in R, with 1000  
355 bootstrap replicates performed to obtain branch support proba-  
356 bilities. Distances between trees were compared using symmetric  
357 Robinson-Foulds distances (Robinson and Foulds 1981) calculated  
358 using the *Phangorn* package in R (Schliep 2011). Cladograms were  
359 plotted using *FigTree* 1.4.2 (<http://tree.bio.ed.ac.uk/software/figtree/>)  
360 with branches transformed to equal lengths for displaying topol-  
361 ogy.

### 362 Data Availability

363 Sequence data used in this study was deposited in the NCBI Short  
364 Read Archive (SRA) under two Bioprojects: PRJNA385842 and  
365 PRJNA385849. The submitted data has been de-multiplexed and  
366 processed through quality control using RedRep.

## 367 RESULTS

### 368 Genotyping-by-sequencing of *NgoMIV-Csp6I* loci in maize

369 *In silico* digestion of the B73 reference genome with *NgoMIV* and  
370 *Csp6I* identified 321,927 sequences that were 64-384 bp long and  
371 were flanked by the recognition site of each enzyme, of which  
372 78,903 were 4PD loci. A 4PD locus in maize ensures that in most  
373 cases the distance between loci within the reference genome is  
374 greater than the average distance within a locus across genomes  
375 (on average, SNPs occur every 28 base pairs in maize [Tenailon  
376 *et al.* 2001]). This was expected to minimize genotyping errors from  
377 ambiguous mapping. A lower PD threshold or a fixed 1 bp differ-  
378 ence threshold may be appropriate and result in the identification  
379 of more SNPs, but we did not consider different thresholds in this  
380 study.

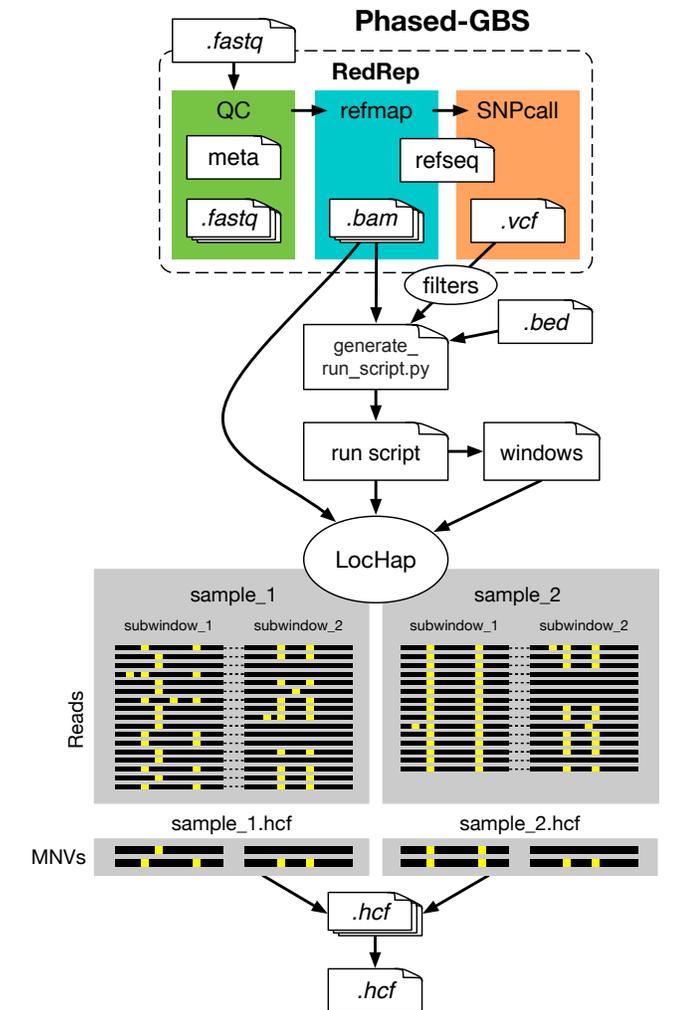
381 For GBS, short-read sequence libraries typically comprise frag-  
 382 ments that fall within a size range of 100 bp. To determine the use  
 383 of targeted size selection to optimize the recovery of fragments  
 384 from non-repetitive loci, we looked at the relative abundance of  
 385 *Ngo*MIV-*Csp*6I loci classified as 4PD versus repetitive ( $\leq 4$ PD) in  
 386 100 bp windows (sliding every 10 bp from 64 bp to 384 bp). The  
 387 median number of 4PD loci among windows was 22,854 (Figure  
 388 S1), but the distribution was somewhat skewed where more 4PD  
 389 loci were present in 100 bp windows of smaller fragment sizes.  
 390 The ratio of 4PD:repetitive loci in 100 bp windows ranged from  
 391  $\sim 0.3$ - $0.4$ , suggesting size selection might be used to maximize se-  
 392 quencing resources by avoiding repetitive sequences. However,  
 393 many restriction endonucleases, including *Ngo*MIV, are sensitive  
 394 to different types of methylation, and SNPs are not expected to be  
 395 uniformly distributed across the genome, such that expectations  
 396 from *in silico* analyses are only a proxy for the numbers of scorable  
 397 loci for a given choice of enzyme. Moreover, the fragments en-  
 398 riched by size selection and PCR are not a uniform sample of the  
 399 underlying distribution of digested fragments and only a fraction  
 400 of those loci will have sufficient read depths for scoring genotypes  
 401 across samples.

402 The genic space (defined here as the gene plus 5 kb flanking  
 403 sequences) comprises  $\sim 27\%$  of the B73 v2 genome (genes alone  
 404 comprise  $\sim 8\%$ ). Unfiltered *Ngo*MIV-*Csp*6I loci from *in silico* diges-  
 405 tion were distributed across the genic and intergenic spaces similar  
 406 to that expected by chance alone (although significantly different,  
 407 coverage of the intergenic space was greater than expected by only  
 408 four percentage points). In total, these unfiltered loci encompass  
 409 3% of the genome and are associated with a majority of the gene  
 410 space, including 69% (27,641) of all maize genes. Filtering *Ngo*MIV-  
 411 *Csp*6I loci that were less than four percent distant from at least one  
 412 alternative site in the genome removed loci associated with 7,472  
 413 genes. The 4PD loci used to score SNPs were enriched within the  
 414 gene space by 32 percentage points (59% for 4PD loci versus an  
 415 expected genome distribution of 27%; Figure 2A) and associated  
 416 with 51% (20,169) of all genes.

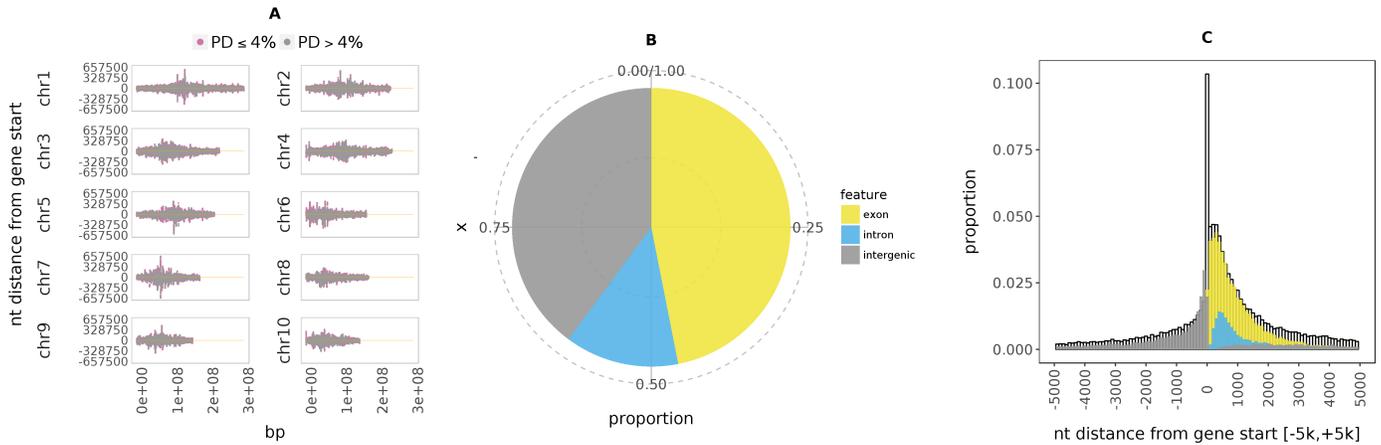
417 A total of 159,994 post-filtered SNPs ( $\geq 3$ X depth at 4PD loci)  
 418 were identified across all samples genotyped in this study (File  
 419 S4; we have also included coordinates for SNPs identified using  
 420 the B73 v3 and B73 v4 reference genomes). These were present  
 421 in 30,239 loci with a median of 4 and mean of 5 SNPs per locus.  
 422 As were the loci themselves, these SNPs were more abundant  
 423 within genes (Figure 2B) and showed enrichment near start codons  
 424 (Figure 2C). The median and mean physical distance between  
 425 loci was 5.7 and 67.9 kb respectively, and the median and mean  
 426 physical distance between SNPs was 14 and 12,854 bp, respectively.  
 427 Previous studies, based on surveys of diverse samples of maize  
 428 have catalogued 1,230 (MaizeSNP50 BeadChip; Ganai *et al.* 2011)  
 429 and 64,662 (HapMap2; Chia *et al.* 2012) of the SNPs found here.

### 430 Empirical evaluation of the accuracy of GBS on heterozygous 431 samples

432 The repeatability and accuracy of GBS was assessed using replicate  
 433 samples and parent-hybrid trios, respectively (Table 1). At 12X  
 434 coverage, the repeatability of genotyping the same DNA was 99.8%  
 435 (based on 265,664 total comparisons). The accuracy of genotyping  
 436 assessed in trios (i.e. observed parent genotypes serving as the  
 437 expectation for predicting an  $F_1$  genotype) was  $>99.6\%$  across all  
 438 monomorphic and polymorphic loci. For expected heterozygotes  
 439 specifically (i.e. only sites that were polymorphic between par-  
 440 ents of a hybrid), accuracy was  $>99.0\%$ . Inspection of incorrectly  
 441 predicted genotypes revealed discrepancies from either putative



**Figure 1 Phased genotyping-by-sequencing.** A reduced representation (RedRep) pipeline is used for SNP typing. The figure shows the basic flow of RedRep, which begins with a *fastq* data file with barcoded sequences. For quality control (QC), the "meta" file contains metadata used for demultiplexing into sample-specific *fastq* files. A reference genome sequence file is used for mapping (refmap) and variant calling (SNPcall). LocHap-GBS is run by editing a generate.py file specifying the location of the *bam* files, the filtered *vcf* file, and a *bed* file of window coordinates to search for haplotypes. A LocHap-GBS run file and windows file are automatically generated. Windows are currently split into subwindows with a maximum of three heterozygous sites within any one individual in the *vcf* file. This situation is depicted for reads across a window that has been delineated into two subwindows where phasing is performed. The dashed connecting line between reads indicates that a contiguous sequence with five SNPs was split into two subwindows. Black-filled bars represent the reference sequence and yellow squares represent SNVs. Given stacks of reads across each subwindow, LocHap-GBS uses a probabilistic model to identify haplotypes in the presence of sequence errors (depicted as one-off instances in the stacks of reads). An *hcf* file is created for each sample, which is then merged into a combined *hcf* file for downstream analysis.



**Figure 2** Genomic distribution of *NgoMIV-Csp6I* loci and SNPs relative to genes in maize. **(A)** Nucleotide distance from each *NgoMIV-Csp6I* *in silico* digested locus to its nearest start codon in the genome (y-axis), relative to the genomic location of the *Csp6I* cut site (x-axis). Two categories of percent distance (PD) are summarized: "unique" loci (greater than 4%; grey points) are overplotted on "repetitive" loci (less than or equal to 4%; red points). **(B)** Proportional distribution of all SNPs discovered in this study with respect to genome annotation categories (~60% fall within genes). **(C)** Distribution of the distance between SNPs and the nearest start site of a gene (shown: 80% of all SNPs which were located within 5 kb of a start codon). Negative values indicate SNPs upstream of the start site or 5'-end of the gene. The distribution of *in silico* digested 4PD loci is plotted as open black bars. Proportions were computed separately for SNPs and *in silico* digested loci, as a function of their respective totals.

442 heterogeneity of the lines or apparent mis-scoring by GATK. There 476  
 443 were one, 12, and zero markers that were monomorphic between 477  
 444 the parents of CML373 X CML341, CML341 X CML277, and Tzi9 478  
 445 X CML258, respectively, which were called as heterozygotes in 479  
 446 the hybrid (putative heterogeneity). There were 26, 52, and 24 480  
 447 polymorphic SNPs between the corresponding parents that were 481  
 448 called as homozygotes in the hybrid. Among these SNPs, there 482  
 449 were 9, 41, and 10, respectively, that had a skewed SNV read-depth 483  
 450 distribution in the hybrid (the minor SNV read depth proportion 484  
 451 was no greater than ~6% of the total read count), despite a median 485  
 452 read-depth of 45X at these loci. For the other 17, 11, and 14 SNPs, 486  
 453 where the minor SNV read depth proportion was relatively high 487  
 454 (minimum ~31%) and there were >10 reads for the minor SNV, the 488  
 455 base-quality score at the variant sites in many of the reads was too 489  
 456 low to be considered in the genotype call. Thus, considering these 490  
 457 latter discrepancies as bonafide inaccuracies in genotyping, along 491  
 458 with the repeatability error rate, the overall accuracy of GBS as 492  
 459 applied in this study was estimated as 99.4%. 493

#### 460 Surveying the same loci across GBS libraries

461 When working with population samples that are larger than the 497  
 462 multiplex size of our GBS design, we found that capturing the 498  
 463 same SNP loci across independently constructed libraries was a 499  
 464 challenge. This is a critical issue for studies where imputation of 500  
 465 missing data is not an option or has unacceptably low accuracy. 501  
 466 Our initial protocol for GBS used a single size-selection step be- 502  
 467 fore PCR amplification and sequencing. Automated size-selection 503  
 468 (BluePippin; Sage Science, Inc.), with the same type of gel cassette 504  
 469 and under the same run settings, yielded differences in the peak 505  
 470 and range of the distributions of size-selected fragments as high 506  
 471 as 75 bp and 41 bp, respectively (Table S1). Introducing a second 507  
 472 size-selection step after PCR minimized these differences to 38 bp 508  
 473 and 36 bp, respectively (Table S1). Not surprisingly, libraries with 509  
 474 similar insert size ranges showed greater correspondence in scored 510  
 475 SNPs than those with dissimilar insert size ranges (Table S2). 511

#### Read-based haplotyping using LocHap-GBS

Variant calling was performed on 276 samples that included four 512  
 test trios (each in duplicate), 23 parents of the maize NAM popula- 513  
 tion, and 234 F<sub>2</sub> individuals derived from the CML373 x CML341 514  
 hybrid (associated with one of the trios). After filtering, there were 515  
 29,706 bi-allelic SNPs (12X read depth) across 9,667 digestion loci. 516  
 In LocHap-GBS, the loci where reads stack up and searches for 517  
 phased MNVs occur, are referred to as windows. Currently, phasing 518  
 in LocHap-GBS is limited to two to three jointly heterozygous 519  
 sites (phasing of more heterozygous sites is under development). 520  
 Because individuals may be heterozygous across more than three 521  
 sites within a window, LocHap-GBS automates the determination 522  
 of subwindows that are delineated based on the nearest sets of 523  
 three (or two) SNPs that are found to be jointly heterozygous in 524  
 any one individual in the *vcf* file. However, subwindows and 525  
 outputted MNVs can be derived from more than the SNP sites that 526  
 define a subwindow. This happens for SNPs within a subwindow 527  
 that are found only in the homozygous state across all individuals 528  
 in the data set, in which case MNVs can exceed a length of three 529  
 nucleotides. As an extreme example, subwindows would corre- 530  
 spond to the original windows for a set of inbred lines where each 531  
 individual is homozygous for every SNP, such that the length of 532  
 MNVs would equal the number of SNP sites within a window. 533  
 Also, the calling routine maintains all of the data in the *vcf* input 534  
 such that some subwindows may capture individual SNPs and the 535  
*hef* output may contain a mixture of MNPs and SNPs. 536

In the 9,667 digestion loci with SNPs, 10,693 subwindows were 537  
 delineated that included 7,749 MNPs and 2,944 SNPs. The number 538  
 of SNPs within MNP subwindows ranged from 2 to 15, with ~90% 539  
 of them having five or fewer SNPs. For MNPs, three subwindows 540  
 contained the observed maximum of 14 MNVs among individ- 541  
 uals, while ~91% of the subwindows had five or fewer MNVs. 542  
 LocHap-GBS led to a 3X reduction in the dimensionality of the 543  
 genotype matrix and a 1.5X increase in the median number of 544  
 variants per locus (an increase from 2 to 3 variants per locus). 545

The trios were used to confirm whether the MNP genotypes

■ **Table 1 Assessment of Typing Accuracy for GBS**

Description	CML373 X CML342	CML341 X CML277	Tzi9 X CML258
A Bi-allelic SNPS	29,706	29,706	29,706
B Not missing in trio	26,216	24,819	21,927
C Consistent across replicates	26,114	24,729	21,878
D Genotyping error: 1-(C/B)	0.004	0.004	0.002
E Not heterozygous in either parent	25,811	24,420	21,555
F Called accurately for SNPs in E	25,784	24,356	21,531
G Genotype accuracy: F/E	0.999	0.997	0.999
H Polymorphic between parents	5,698	5,004	5,373
I Called accurately for loci in H	5,672	4,952	5,349
J Heterozygote genotype accuracy: I/H	0.995	0.990	0.996
K Not called as heterozygote for H: H-I	26	52	24
L Low minor allele depth in K	9	41	10
M Approximately equal allele depth in K	17	11	14
N Adjusted genotype accuracy: G-D	0.995	0.994	0.997

scored by LocHap-GBS were valid. For each trio, subwindows were excluded if they contained missing data, had inconsistent calls between replicate samples, or contained heterozygous SNP genotypes in the parents. The MNVs expected to be observed in the hybrid were determined directly from the filtered SNP data on the parental lines (not by LocHap-GBS). These were then compared to MNVs called by LocHap-GBS. Among 25,404 MNVs recorded for the parents of all trios (ignoring whether MNVs were shared between trios), six were scored differently by LocHap-GBS in the hybrids. All of these were associated with the inaccurately scored genotypes noted previously.

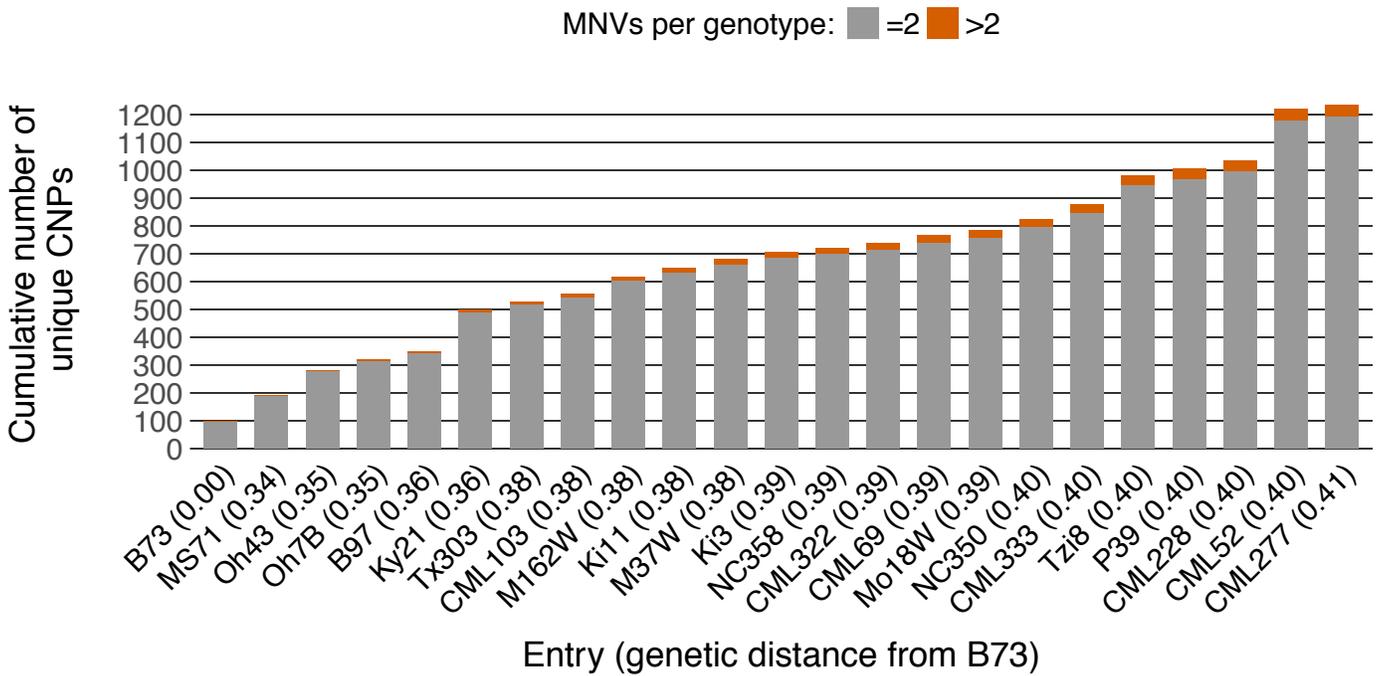
When evaluating unique haplotype numbers for each subwindow we noticed some MNP loci in the hybrids that had more than two MNVs, which is not possible for maize given that it is diploid. Inspection of the raw sequence data at these loci revealed no barcode swapping errors. Examining the transmission of MNVs in the trios and F<sub>2</sub> population allowed us to rule out DNA cross-contamination and conclude that most, if not all of these corresponded to divergent paralogs in the sequenced samples that collapsed onto a single locus from the B73 reference genome; we refer to these as CNPs. For instance, a CML373 X CML341 F<sub>1</sub> sample had 21 loci with three or four MNVs that were consistent in replicate samples. For all of these loci, one or both of the inbred parents had more than one MNV (this was determined by maintaining heterozygous SNP genotypes in the inbred parents, which were filtered when estimating accuracy above) which matched each of the MNVs found in the hybrid. These same MNVs were found in the F<sub>2</sub> population, and tests of co-segregation indicated that 11 of them were genetically linked, but some of these loci included a mixture of linked and unlinked MNVs (data not shown).

Under the assumption that loci with greater than one MNV in an inbred line represent CNPs, among the 23 inbred parents of the NAM population there was an average of 155 (2.0%) total CNPs per individual (maximum of 305 CNPs for CML52). The cumulative number of unique CNPs was recorded with the addition of each inbred line inserted in order of their estimated genetic distance

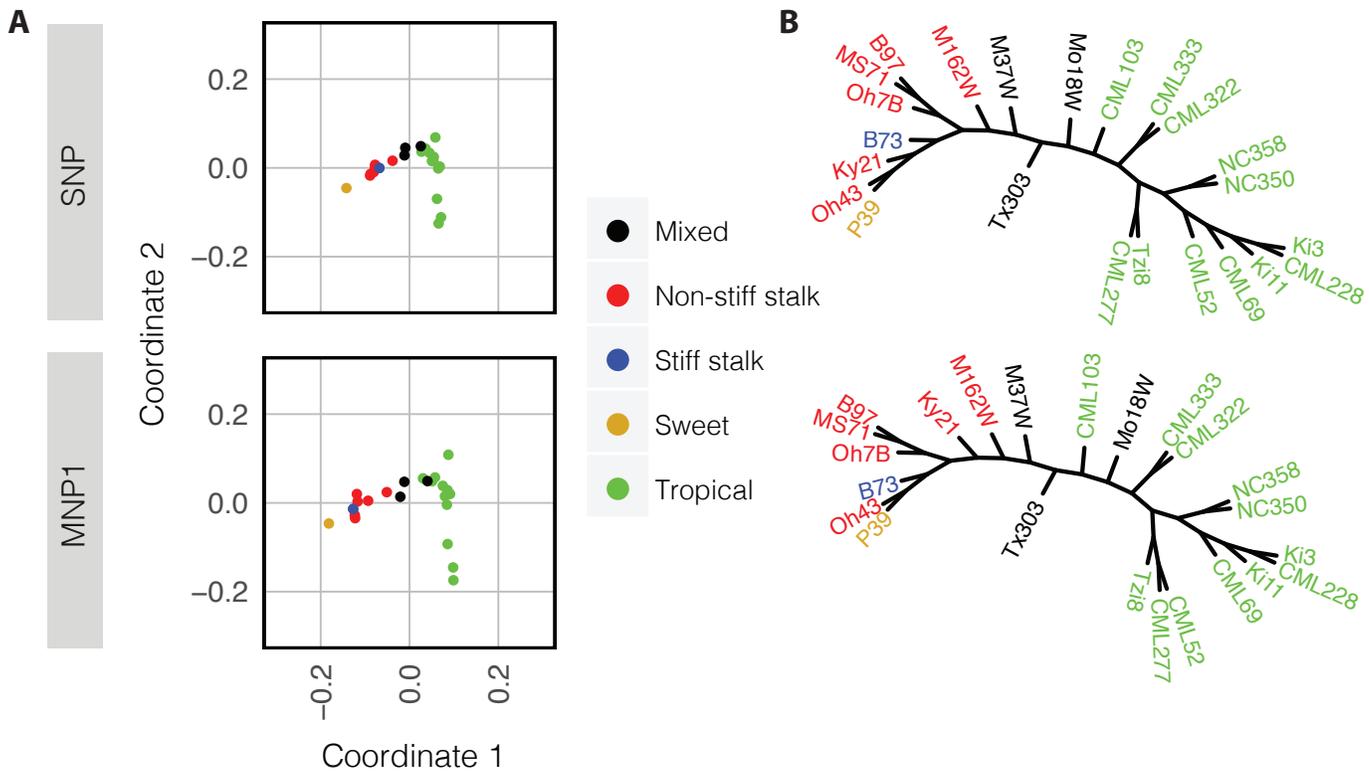
from B73 (Figure 3). There was a median increase of 30 unique CNP-associated loci with exactly two MNVs, culminating in a total of 1,195 such loci among all of the lines (Figure 3). This is an upper-bound estimate of the number of CNPs, since CNPs are conflated with heterozygous loci in inbred lines. For a lower-bound estimate, loci with more than two MNVs per line were recorded, which showed an average of 2 (0.03%) total CNPs per line (maximum of 8), a median increase of 2 CNPs with each additional line added to the dataset, and a cumulative total of 41 unique CNPs. Taken together, the estimated proportion of CNPs typed in this study among the 23 parental lines of the NAM population lies between 0.5% and 15.4% (Figure 3). Excluding the reference line B73, the correlation between genetic distance and number of unique CNPs per additional line was not significantly different from zero for both classifications of CNPs.

### Genetic diversity based on SNPs versus MNPs

Results from the analysis of genetic diversity among parents of the NAM population were compared using SNP and MNP datasets. Excluding CNP-associated genotypes (MNP2 and MNP3 datasets) from the MNP1 dataset showed no significant difference in the distribution of pairwise genetic distances and all three datasets were almost perfectly correlated ( $r > 0.999$ ; Mantel test:  $p = 0.001$ ). Therefore, the following results are reported only for the SNP and MNP1 datasets. Although the distance matrices were significantly correlated ( $r = 0.890$ ; Mantel test:  $p = 0.001$ ), the average pairwise distance for SNPs was lower than that for MNPs (0.28 and 0.37 respectively) and the range in the distances was smaller for SNPs (0.12) than for MNPs (0.20). Although the mean bootstrap probabilities were effectively identical for the consensus trees based on SNPs (79.3%) and MNPs (78.9%), Robinson-Foulds distance between those trees was 12 (maximum possible distance was 40). Consequently, multidimensional scaling with MNPs produced greater separation among the lines (Figure 4A) and resulted in differences in their topology (Figure 4B; e.g., c.f. CML103, Ky21 and CML52).



**Figure 3 CNP-associated loci in the NAM inbred parents.** From left-to-right, each bar indicates the cumulative number of unique CNPs (loci containing more than one MNV) in each individual compared to its preceding set of lines, with the lines ordered by their MNP1-based genetic distance from B73. The grey portion of the bar is the number of loci that had exactly two MNVs and the orange portion of the bar is the number of loci that had more than two MNVs.



**Figure 4 Analysis of genetic diversity with SNP and MNP data.** (A) MDS plots based on shared allele distance for SNP (top) and MNP1 (bottom) data on 23 maize inbred lines. (B) Corresponding Neighbor-Joining trees. Colors represent previously assigned population structure (Liu *et al.* 2003).

584 Many implementations of GBS are not optimal for genotyping of  
 585 heterozygous individuals or populations with unknown parentage.  
 586 Imputation has proven successful on biparental families of recom-  
 587 binant inbred lines, but falls short when applied to populations that  
 588 have low linkage disequilibrium (the median imputation accuracy  
 589 was shown to be between 15-80% for markers with a  $r^2 < 0.3$ ; He  
 590 *et al.* 2015). The few approaches that have been developed for gen-  
 591 otyping heterozygous populations use population-specific filters  
 592 and depend on relatedness among individuals to score genotypes  
 593 or require targeted sequencing to obtain read depths sufficient for  
 594 calling heterozygotes (Uitdewilligen *et al.* 2013; Gardner *et al.* 2014;  
 595 Barba *et al.* 2014; Hyma *et al.* 2015). HetMappS (Hyma *et al.* 2015)  
 596 was developed specifically for dense genetic maps, and markers  
 597 are filtered based on expected genotype ratios for pseudo-testcross-  
 598 markers. This method was successful at increasing marker density,  
 599 but was designed for use in F<sub>1</sub> populations only. Here, we have  
 600 presented a GBS protocol and bioinformatic pipeline that, at least  
 601 for maize, produces highly accurate genotype data on heterozy-  
 602 gous individuals without requiring information on parentage or  
 603 family structure nor imputation to obtain high call rates on the  
 604 typed SNPs or MNPs.

605 As discussed by Peterson *et al.* (2012), there is a balancing act  
 606 in deciding which enzymes and size-selection windows should  
 607 be used when implementing GBS. In maize, which has been exten-  
 608 sively explored for SNPs, where the enzyme *ApeKI* has been  
 609 used routinely for GBS, the *NgomIV-Csp6I* enzyme combination  
 610 provided an effective means for accurate genotyping that led to  
 611 the discovery of many new SNPs. Assuming there is a 50% chance  
 612 of sequencing each allele at a bi-allelic marker, a minimum read  
 613 depth of 12X predicts that the binomial probability of sequencing  
 614 each allele at least twice is >99.6%. In this study, accuracy was  
 615 estimated experimentally using parent-hybrid trios to be >99.4%  
 616 (Table 1), which fit closely to the expected sampling probability.  
 617 Compared to an existing catalogue of 52,340,265 SNPs (Ganal  
 618 *et al.* 2011; Chia *et al.* 2012) identified among diverse individuals of  
 619 maize that included all but four of the ones used here (excluding:  
 620 CML10, CML258, CML373, Tzi9), this study discovered approxi-  
 621 mately 100,000 new SNPs which comprised nearly two-thirds of  
 622 the total SNPs found. These were enriched near genes, which  
 623 was in part due to the distribution of 4PD loci, but may also be  
 624 attributed to patterns in methylation: *NgomIV* is averse to the CpG  
 625 methylation which is predominant in the intergenic repeat space  
 626 of the maize genome (Antequera and Bird, 1988). Enrichment of  
 627 the typed loci around the start codon of genes (Figure 2C) might  
 628 be because the starts of monocot genes are GC-rich (Glémin *et al.*  
 629 2014) and *NgomIV* recognizes a GC-rich recognition site.

630 With the decline in cost for HTS, there is growing opportunity to  
 631 apply GBS to thousands of samples for genetic studies. Sequencing  
 632 the same loci within and across separate GBS libraries is required  
 633 for this, but whether the same loci are sequenced has not been  
 634 examined much in the literature (we are aware of one exception:  
 635 DaCosta and Sorenson 2014); perhaps because GBS studies tend  
 636 to rely on imputation of missing data. It has been reasoned that  
 637 size selection would play an important role in sequencing shared  
 638 loci (Ali *et al.* 2016). Peterson *et al.* (2012) showed that narrower  
 639 size-selection windows increased the sequencing depth of shared  
 640 loci in a library. As a cautionary note, some size-selection win-  
 641 dows contain an abundance of repetitive loci (Figure S1). Here,  
 642 we addressed the issue of sampling shared loci across different  
 643 libraries, reaching the conclusion that two size-selection steps were  
 644 needed to maximize the number of shared loci sequenced. Im-

645 portantly, our results led us to the realization that the standard  
 646 Y-adapters used for HTS libraries may form structures that migrate  
 647 slower than fully dsDNA in dye-free agarose gels. This can lead  
 648 to inconsistent size selection when using automated size-selection  
 649 instrumentation (Sadaf Hoda; Sage Science, pers. comm.). Because  
 650 of the inconsistency in size selection of libraries with Y-adaptors,  
 651 we had to reduce the pre-PCR size-selection window and introduce  
 652 a post-PCR size-selection on the dsDNA created by PCR ampli-  
 653 fication. This minimized the variation in size-selection between  
 654 libraries and maximized the recovery of shared GBS loci.

655 This study expands the usefulness of GBS data, allowing for  
 656 phased genotyping of MNPs, which enhances the information  
 657 available for genetic studies. For example, founder haplotypes  
 658 that are resolved for multi-parental populations can be used to  
 659 impute missing data in the progeny (Gatti *et al.* 2014) or to initial-  
 660 ize imputation (Davies *et al.* 2016). Also, the use of haplotypes  
 661 may prove useful in uncovering regions associated with adaptive  
 662 traits in both model and non-model species (Lorenz *et al.* 2010).  
 663 The extension of LocHap for GBS was developed as a new tool  
 664 for phased genotyping using standard file formats produced by  
 665 most mapping and variant calling software. LocHap-GBS may be  
 666 advantageous over other current tools capable of scoring haplo-  
 667 types, such as STACKS and Haplotag, in that statistically vetted  
 668 MNVs can be scored by LocHap-GBS on individuals without popu-  
 669 lation filters only definable for specific types of samples. However,  
 670 LocHap-GBS may not be readily integrated by tools that use cus-  
 671 tomized data formats, though the algorithm developed for local  
 672 haplotyping could (Sengupta *et al.* 2015). Finally, the software  
 673 GATK used to score SNPs can also generate phased genotypes  
 674 (McKenna *et al.* 2010). RedRep (Figure 1) was developed using  
 675 the GATK UnifiedGenotyper, which does not have MNP phasing  
 676 functionality (GATK HaplotypeCaller does). We have not experi-  
 677 mented with the current phasing tool of GATK, but we note some  
 678 unique features of LocHap-GBS. It uses a distinct algorithm for  
 679 identifying MNVs, can be run independently with standard input  
 680 files, reports phased genotypes in a different format (which we  
 681 consider more appropriate for direct downstream analysis), and  
 682 can score more than two haplotypes at a locus—this led to our  
 683 finding of CNPs in maize.

684 Phased genotype data, while decreasing the marker dimension-  
 685 ality of the genotype matrix, was more informative than SNP data.  
 686 Previously characterized lines of maize were more differentiated  
 687 and showed some differences in relatedness when analyzed using  
 688 MNPs compared to SNPs (Figure 4). In addition, phased-GBS  
 689 revealed CNPs that were validated by genetic transmission yet  
 690 undetectable in the underlying SNP data. These can be explained  
 691 as divergent paralogs within the sample genomes that are present  
 692 as only one copy in the reference genome. Because the SNP data  
 693 was filtered to examine only 4PD loci, it seems unlikely that the  
 694 CNPs were a result of mapping to sequences that have been col-  
 695 lapsed in the assembly of the B73 reference genome, but this is a  
 696 possibility that could result in reproducibly inaccurate CNPs. Nev-  
 697 ertheless, having included B73 in our set of sequenced individuals,  
 698 which is the same line used to sequence the reference genome, we  
 699 determined that our B73 sample harbored no GBS loci with more  
 700 than two MNVs, while other inbred lines did (Figure 3). There  
 701 were, however, nearly 100 GBS loci that were heterozygous in the  
 702 B73 sample. These may be attributed to residual heterozygosity  
 703 in the line, errors in the assembly, or CNPs that had arisen during  
 704 seed increases of the stock. Inspecting the genomic distribution of  
 705 these loci showed no pattern of clustering as one might expect for  
 706 residual heterozygosity, suggesting one or both of the latter two

707 explanations.

708 Although our study was not designed for CNP discovery, it led  
709 to the identification of at least 41 CNP-associated loci (potentially  
710 around one thousand) among 23 inbred lines of maize. The average  
711 read-depth at CNP-associated loci was 77, which was about the  
712 average read-depth of non-CNP associated loci. However, the read-  
713 depth range for CNP-associated loci was 1,480 (max: 1,492 reads)  
714 compared to 755 for non-CNP associated loci (max: 767 reads),  
715 which is expected if there are multiple paralogous loci mapping  
716 to the same place in the reference genome. Moreover, despite an  
717 average depth of sequence coverage that gives a high likelihood for  
718 detecting at least 10 CNVs, no more than four CNVs at a locus were  
719 found within an individual. There was no relationship between  
720 genome-wide estimates of genetic distance and the number of  
721 unique CNP-associated loci (Figure 3), suggesting the evolution  
722 of these putative CNPs may be different from SNPs and MNPs.  
723 This study demonstrates the potential of extending GBS for phased  
724 genotyping. As applications of GBS expand beyond bi-parental  
725 mapping populations, we foresee numerous benefits to typing  
726 MNPs by phased-GBS.

## 727 ACKNOWLEDGEMENTS

728 This project was supported by the Agriculture and Food Research  
729 Initiative Competitive Grant No. 2011-67003-30342 from the USDA  
730 National Institute of Food and Agriculture (Agriculture and Natu-  
731 ral Resources Science for Climate Variability and Change Program).  
732 Other support from the University of Delaware Bioinformatics  
733 Core Facility, including use of the BIOMIX computational cluster,  
734 was made possible by Delaware INBRE (NIH P20 GM103446). We  
735 thank Dr. Karol Miaskiewicz at the Delaware Biotechnology Insti-  
736 tute for assistance using BIOMIX. We thank Dr. Patricia Klein at  
737 TAMU for suggesting we consider *NgoMIV*.

## 738 LITERATURE CITED

739 Ali, O. A., S. M. O'Rourke, S. J. Amish, M. H. Meek, G. Luikart, *et al.*,  
740 2016 RAD Capture (Rapture): flexible and efficient sequence-  
741 based genotyping. *Genetics* **202**: 389–400.  
742 Andolfatto, P., D. Davison, D. Erezylmaz, T. T. Hu, J. Mast, *et al.*,  
743 2011 Multiplexed shotgun genotyping for rapid and efficient  
744 genetic mapping. *Genome Res* **21**: 610–617.  
745 Andrews, K. R., J. M. Good, M. R. Miller, G. Luikart, and P. A.  
746 Hohenlohe, 2016 Harnessing the power of RADseq for ecological  
747 and evolutionary genomics. *Nat Rev Genet* **17**: 81–92.  
748 Baird, N. A., P. D. Etter, T. S. Atwood, M. C. Currey, A. L. Shiver,  
749 *et al.*, 2008 Rapid SNP discovery and genetic mapping using  
750 sequenced RAD markers. *PLoS One* **3**.  
751 Barba, P., L. Cadle-Davidson, J. Harriman, J. C. Glaubitz, S. Brooks,  
752 *et al.*, 2014 Grapevine powdery mildew resistance and suscepti-  
753 bility loci identified on a high-resolution SNP map. *Theor Appl*  
754 *Genet* **127**: 73–84.  
755 Beissinger, T. M., C. N. Hirsch, R. S. Sekhon, J. M. Foerster, J. M.  
756 Johnson, *et al.*, 2013 Marker density and read depth for genotyp-  
757 ing populations using genotyping-by-sequencing. *Genetics* **193**:  
758 1073–1081.  
759 Bowcock, A. M., A. Ruiz-Linares, J. Tomfohrde, E. Minch, J. R.  
760 Kidd, *et al.*, 1994 High resolution of human evolutionary trees  
761 with polymorphic microsatellites. *Nature* **368**: 455–457.  
762 Catchen, J., P. A. Hohenlohe, S. Bassham, A. Amores, and W. A.  
763 Cresko, 2013 Stacks: an analysis tool set for population genomics.  
764 *Mol Ecol* **22**: 3124–3140.

765 Chia, J. M., C. Song, P. J. Bradbury, D. Costich, N. de Leon, *et al.*,  
766 2012 Maize HapMap2 identifies extant variation from a genome  
767 in flux. *Nat Genet* **44**: 803–807.  
768 DaCosta, J. M. and M. D. Sorenson, 2014 Amplification biases and  
769 consistent recovery of loci in a double-digest RAD-seq protocol.  
770 *PLoS One* **9**.  
771 Davies, R. W., J. Flint, S. Myers, and R. Mott, 2016 Rapid genotype  
772 imputation from sequence without reference panels. *Nature*  
773 *Genetics* **48**: 965–969.  
774 Elshire, R. J., J. C. Glaubitz, Q. Sun, J. A. Poland, K. Kawamoto,  
775 *et al.*, 2011 A robust, simple genotyping-by-sequencing (GBS)  
776 approach for high diversity species. *PLoS One* **6**.  
777 Ferdosi, M. H., J. Henshall, and B. Tier, 2016 Study of the optimum  
778 haplotype length to build genomic relationship matrices. *Genet*  
779 *Sel Evol* **48**.  
780 Fu, Y.-B. and G. W. Peterson, 2011 Genetic diversity analysis with  
781 454 pyrosequencing and genomic reduction confirmed the east-  
782 ern and western division in the cultivated barley gene pool.  
783 *Plant Genome* **4**: 226–237.  
784 Ganai, M. W., G. Durstewitz, A. Polley, A. Bérard, E. S. Buckler,  
785 *et al.*, 2011 A large maize (*Zea mays* L.) SNP genotyping array:  
786 development and germplasm genotyping, and genetic mapping  
787 to compare with the B73 reference genome. *PLoS One* **6**.  
788 Gardner, K. M., P. J. Brown, T. F. Cooke, S. Cann, F. Costa, *et al.*,  
789 2014 Fast and cost-effective genetic mapping in apple using  
790 next-generation sequencing. *G3* **4**: 1681–1687.  
791 Gascuel, O., 1997 BIONJ: an improved version of the NJ algorithm  
792 based on a simple model of sequence data. *Mol Biol Evol* **14**:  
793 685–695.  
794 Gatti, D. M., K. L. Svenson, A. Shabalina, L.-Y. Wu, W. Valdar, *et al.*,  
795 2014 Quantitative trait locus mapping methods for Diversity  
796 Outbred Mice. *G3* **4**: 1623–1633.  
797 He, S., Y. Zhao, M. F. Mette, R. Bothe, E. Ebmeyer, T. F. Sharbel,  
798 J. C. Reif, and Y. Jiang, 2015 Prospects and limits of marker im-  
799 putation in quantitative genetic studies in european elite wheat  
800 (*triticum aestivum* l.). *BMC genomics* **16**: 168.  
801 Heffelfinger, C., C. A. Fragoso, M. A. Moreno, J. D. Overton,  
802 J. P. Mottinger, *et al.*, 2014 Flexible and scalable genotyping-  
803 by-sequencing strategies for population studies. *BMC Genomics*  
804 **15**: 979.  
805 Hyma, K. E., P. Barba, M. Wang, J. P. Londo, C. B. Acharya, *et al.*,  
806 2015 Heterozygous mapping strategy (HetMappS) for high res-  
807 olution genotyping-by-sequencing markers: A case study in  
808 grapevine. *PLoS One* **10**.  
809 Jombart, T., 2008 adegenet: a R package for the multivariate analy-  
810 sis of genetic markers. *Bioinformatics* **24**: 1403–1405.  
811 Kruskal, J. B., 1964 Multidimensional scaling by optimizing good-  
812 ness of fit to a nonmetric hypothesis. *Psychometrika* **29**: 1–27.  
813 Li, H., 2013 Aligning sequence reads, clone sequences and assem-  
814 bly contigs with BWA-MEM. *arXiv:1303.3997* .  
815 Liu, K., M. Goodman, S. Muse, J. S. Smith, E. Buckler, *et al.*, 2003  
816 Genetic structure and diversity among maize inbred lines as  
817 inferred from DNA microsatellites. *Genetics* **165**: 2117–2128.  
818 Lorenz, A. J., M. T. Hamblin, and J.-L. Jannink, 2010 Performance of  
819 single nucleotide polymorphisms versus haplotypes for genome-  
820 wide association analysis in barley. *PIOS ONE* **5**.  
821 Lu, Y., T. Shah, Z. Hao, S. Taba, S. Zhang, *et al.*, 2011 Comparative  
822 SNP and haplotype analysis reveals a higher genetic diversity  
823 and rapider LD decay in tropical than temperate germplasm in  
824 maize. *PLoS One* **6**.  
825 Lund, A. H., M. Duch, and F. S. Pedersen, 1996 Increased cloning  
826 efficiency by temperature-cycle ligation. *Nucleic Acids Res* **24**:

800–801.

827 Mantel, N., 1967 The detection of disease clustering and a general-  
828 ized regression approach. *Cancer research* **27**: 209–220.  
829

830 Martin, M., 2011 Cutadapt removes adapter sequences from high-  
831 throughputs sequencing reads. *EMBnetjournal* **17**: 10–12.

832 McKenna, A., M. Hanna, E. Banks, A. Sivachenko, K. Cibul-  
833 skis, *et al.*, 2010 The Genome Analysis Toolkit: a MapReduce  
834 framework for analyzing next-generation DNA sequencing data.  
835 *Genome Res* **20**: 1297–1303.

836 McMullen, M. D., S. Kresovich, H. S. Villeda, P. Bradbury, H. Li,  
837 *et al.*, 2009 Genetic properties of the maize nested association  
838 mapping population. *Science* **325**: 737–740.

839 Neph, S., M. S. Kuehn, A. P. Reynolds, E. Haugen, R. E. Thur-  
840 man, *et al.*, 2012 BEDOPS: high-performance genomic feature  
841 operations. *Bioinformatics* **28**: 1919–1920.

842 Paradis, E., J. Claude, and K. Strimmer, 2004 APE: analyses of  
843 phylogenetics and evolution in R language. *Bioinformatics* **20**:  
844 289–290.

845 Peterson, B. K., J. N. Weber, E. H. Kay, H. S. Fisher, and H. E.  
846 Hoekstra, 2012 Double digest RADseq: an inexpensive method  
847 for de novo SNP discovery and genotyping in model and non-  
848 model species. *PLoS One* **7**.

849 Poland, J. A., P. J. Brown, M. E. Sorrells, and J.-L. Jannink, 2012  
850 Development of high-density genetic maps for barley and wheat  
851 using a novel two-enzyme genotyping-by-sequencing approach.  
852 *PLoS One* **7**.

853 Robinson, D. F. and L. R. Foulds, 1981 Comparison of phylogenetic  
854 trees. *Mathematical Biosci* **53**: 131–147.

855 Rocher, S., M. Jean, Y. Castonguay, and F. Belzile, 2015 Validation of  
856 genotyping-by-sequencing analysis in populations of tetraploid  
857 alfalfa by 454 sequencing. *PLoS One* **10**: 1–18.

858 Schliep, K., 2011 phangorn: phylogenetic analysis in r. *Bioinfor-*  
859 *matics* **27**: 592–593.

860 Sengupta, S., K. Gulukota, Y. Zhu, C. Ober, K. Naughton, *et al.*,  
861 2015 Ultra-fast local-haplotype variant calling using paired-end  
862 DNA-sequencing data reveals somatic mosaicism in tumor and  
863 normal blood samples. *Nucleic Acids Res* **44**.

864 Sonah, H., M. Bastien, E. Iquiria, A. Tardivel, G. Légaré, *et al.*,  
865 2013 An improved genotyping by sequencing (GBS) approach  
866 offering increased versatility and efficiency of SNP discovery  
867 and genotyping. *PLoS One* **8**.

868 Tenaillon, M. I., M. C. Sawkins, A. D. Long, R. L. Gaut, J. F. Doebl-  
869 ley, *et al.*, 2001 Patterns of DNA sequence polymorphism along  
870 chromosome 1 of maize (*Zea mays* ssp. *mays* L.). *P Natl Acad*  
871 *Sci* **98**: 9161–9166.

872 Tinker, N. A., W. A. Bekele, and J. Hattori, 2016 Haplotag: software  
873 for haplotype-based genotyping-by-sequencing analysis. *G3* **6**:  
874 857–863.

875 Torkamaneh, D., J. Laroche, and F. Belzile, 2016 Genome-wide SNP  
876 calling from genotyping by sequencing (GBS) data: A compari-  
877 son of seven pipelines and two sequencing technologies. *PLoS*  
878 *One* **11**.

879 Uitdewilligen, J. G. A. M. L., A. M. A. Wolters, B. B. D’hoop, T. J. A.  
880 Borm, R. G. F. Visser, *et al.*, 2013 A next-generation sequencing  
881 method for genotyping-by-sequencing of highly heterozygous  
882 autotetraploid potato. *PLoS One* **8**.