

Marker imputation efficiency for genotyping-by-sequencing data in rice (*Oryza sativa*) and alfalfa (*Medicago sativa*)

Nelson Nazzicari · Filippo Biscarini ·
Paolo Cozzi · E. Charles Brummer ·
Paolo Annicchiarico

Received: 30 June 2015 / Accepted: 13 May 2016 / Published online: 25 May 2016
© Springer Science+Business Media Dordrecht 2016

Abstract Genotyping-by-sequencing (GBS) is a rapid and cost-effective genome-wide genotyping technique applicable whether a reference genome is available or not. Due to the cost-coverage trade-off, however, GBS typically produces large amounts of missing marker genotypes, whose imputation becomes therefore both challenging and critical for later analyses. In this work, the performance of four general imputation methods (K-nearest neighbors,

Random Forest, singular value decomposition, and mean value) and two genotype-specific methods (“Beagle” and FILLIN) was measured on GBS data from alfalfa (*Medicago sativa* L., autotetraploid, heterozygous, without reference genome) and rice (*Oryza sativa* L., diploid, 100 % homozygous, with reference genome). Alfalfa SNP were aligned on the genome of the closely related species *Medicago truncatula* L.. Benchmarks consisted in progressive data filtering for marker call rate (up to 70 %) and increasing proportions (up to 20 %) of known genotypes masked for imputation. The relative performance was measured as the total proportion of correctly imputed genotypes, globally and within each genotype class (two homozygotes in rice, two homozygotes and one heterozygote in alfalfa). We found that imputation accuracy was robust to increasing missing rates, and consistently higher in rice than in alfalfa. Accuracy was as high as 90–100 % for the major (most frequent) homozygous genotype, but dropped to 80–90 % (rice) and below 30 % (alfalfa) in the minor homozygous genotype. Beagle was the best performing method, both accuracy- and time-wise, in rice. In alfalfa, KNNI and RFI gave the highest accuracies, but KNNI was much faster.

Electronic supplementary material The online version of this article (doi:[10.1007/s11032-016-0490-y](https://doi.org/10.1007/s11032-016-0490-y)) contains supplementary material, which is available to authorized users.

Nelson Nazzicari and Filippo Biscarini contributed equally to the work.

N. Nazzicari (✉)
Council for Agricultural Research and Economics
(CREA) Research Centre for Fodder Crops and Dairy
Productions, Lodi, Italy
e-mail: nelson.nazzicari@crea.gov.it

F. Biscarini · P. Cozzi · P. Annicchiarico
Dipartimento di Bioinformatica, Fondazione Parco
Tecnologico Padano, Lodi, Italy

E. C. Brummer
Plant Sciences Department, University of California,
Davis, CA, USA

Keywords SNP · Genotyping by sequencing (GBS) · K-nearest neighbors imputation (KNNI) · Random Forest imputation (RFI) · Singular value decomposition imputation (SVDI) · Beagle · FILLIN · Alfalfa · Rice · Imputation · Reference genome

Introduction

Imputation of missing alleles and genotypes is a preliminary step for a wide range of genetic analyses. In fact most models and software for population genetics, genomic selection (GS) and genome-wide association studies (GWAS) can not easily handle missing data and require complete datasets (e.g., Hayes et al. 2009; Aulchenko et al. 2007; Endelman 2011; Pérez and de los Campos 2014). SNP array represents one of the leading genotyping technologies and produce datasets that, after low call rate filtering, usually still contain a small proportion of uncalled genotypes (e.g., <5 %) randomly distributed along the genome. Genotyping-by-sequencing (GBS) is a relatively recent technique which is considered a viable alternative to SNP array-based genotyping to produce SNP genotype data (Elshire et al. 2011). GBS is platform independent and is conveniently used in species that lack commercial SNP chips or even lack a reference genome sequence. GBS data typically present a much larger proportion of sporadic missing genotypes: e.g., ~52 % and ~58 % average in two maize experimental populations (Crossa et al. 2013), or “up to 80 % missing data per marker” in wheat (Poland et al. 2012). This is due to both intrinsic properties of the technology and to its application to species without a mature genome assembly (Glaubitz et al. 2014).

There are several methods that are routinely used for the imputation of missing genotypes. Some rely solely on genotypic information, some make use also of genealogies, most of them are based on reconstructing haplotypes to be used in predictive models (Nicolazzi et al. 2015). Most methods for genotype imputation have been applied to SNP array data, and typically yield very high imputation accuracy. For instance, >95 % correctly imputed genotypes were reported in maize Hickey et al. (2012) and cattle VanRaden et al. (2011). The same imputation methods can also be applied to GBS data (see Huang et al. (2014); Swarts et al. (2014) for applications in rice and maize). When a reference genome is available, SNP loci can be aligned against its sequence and are therefore ordered, thus allowing the exploitation of local linkage disequilibrium (LD). However, GBS data may be generated also for species lacking a reference genome assembly. In this case the GBS output is a series of “floating” loci not linked to any chromosome: Alignment is not possible and SNPs are

therefore considered unordered. Unordered data add additional complexity to genotype imputation. Exploitation of linkage disequilibrium and haplotype reconstruction are less straightforward, consecutive loci may lie on different chromosomes or scaffolds, and overall data are less homogeneous. Rutkoski et al. (2013) considered the application of general data imputation methods to wheat, maize and barley populations, obtaining best case scenarios accuracies as high as 0.84–0.94 (measured as R^2 between true and imputed genotypes). Imputation of missing GBS data without alignment to the reference genome has been reported in alfalfa (*Medicago sativa* L., Rocher et al. 2015) and red raspberry (*Rubus idaeus* L., Ward et al. 2013). These studies, however, do not focus on measuring imputation accuracy and do not compare alternative imputation methods.

Compared to SNP array data, GBS is more challenging: many more missing genotypes, high variability between runs, intrinsically noisy data, reads from different loci that can overlap. For all these reasons, imputation of missing genotypes with GBS data is still a relatively immature technique, not yet amenable to produce highly standardized and repeatable results. There is therefore scientific and practical interest in gaining further insights into how genotype imputation with GBS data works, and in developing the best methods and strategies to accurately and efficiently impute such missing data. This would be relevant not only for the scientific community but also for the breeding industry, because of its direct consequences on genomic applications.

In this paper, we imputed missing genotypes from GBS data in two agronomically important crop species: the cereal crop rice (*Oryza sativa* L.) whose reference genome has been assembled in International Rice Genome Sequencing Project (2005), and updated in Kawahara et al. (2013), and the forage crop alfalfa (*M. sativa* L.), for which a reference genome is not available yet. Alfalfa genotypes were mapped on the close relative *Medicago truncatula* L., a model species for legumes, whose genome is available (Young et al. 2011). The species were chosen as representative of very different scenarios, alfalfa being autotetraploid with high heterozygosity and rice being diploid with essentially 100 % homozygosity. The present study thus encompasses a wide span of real application cases.

We compared the imputation accuracy of four general imputation methods (mean value imputation, K-nearest neighbor imputation, singular value decomposition imputation and Random Forest imputation) with the performance of two methods specific for genotype imputation (localized haplotype clustering imputation, implemented in the software package Beagle, and the method of based on haplotype reconstruction around recombination sites implemented in the FILLIN algorithm as part of the Tassel suite). The four general algorithms were chosen as well-known imputation strategies implemented in several freely available software libraries. The two genotype-specific methods represent the state of the art for genotypes lacking pedigree information.

The relative efficiency of the different imputation methods was assessed in terms of accuracy and computation time. Accuracy was measured as fraction of correctly imputed genotypes, both as a total and within genotype classes (major/minor homozygotes, and heterozygotes, when present).

Materials and methods

Plant material

Samples from two autotetraploid alfalfa (*M. sativa* L.) and one diploid rice (*O. sativa* L.) populations were available (see Table 1).

Alfalfa populations included elite germplasms from the Po Valley (Alfalfa-PV) and Mediterranean-climate environments (Alfalfa-Med). For Alfalfa-PV, 124 parent genotypes were chosen by stratified mass selection for dry matter yield over three harvests. For Alfalfa-Med, 154 parent genotypes were derived from two cycles of free intercrossing among three outstanding populations in a multi-environment study. For further details see Annicchiarico et al. (2015).

The Rice dataset included 437 plants belonging to 391 accessions (46 duplicates) from the Rice Germplasm Collection maintained at CRA-Rice Research Unit (Vercelli, Italy). The sampled collection comprised accessions from the five main subpopulations of *O. sativa* (274 temperate japonica, 108 tropical japonica, 28 indica, 16 aus and 11 aromatic). All accessions were purified through single seed descent, and were genotypically essentially 100 % homozygous.

Genotyping by sequencing

The 715 alfalfa (Alfalfa-Med and Alfalfa-PV) and rice (Rice) samples were genotyped using the genotyping-by-sequencing (GBS) approach (Elshire et al. 2011). Different GBS protocols were used to genotype alfalfa and rice, since genotyping was carried out under different projects.

In alfalfa, DNA was isolated from fresh leaf tissues by the Wizard® Genomic DNA Purification Kit (Promega, A1125) and quantified with a Quant-iT PicoGreen dsDNA assay kit (Life Technologies, P7589). Two libraries were constructed for the two populations, where 100 ng of each DNA was digested with ApeKI (NEB, R0643L) and then ligated to a unique barcoded adapter and a common adapter. For the reference population Alfalfa-Med, 5 nM each of the primers and NEB 2X Taq Master Mix (NEB Cat # M0270S) were included in the PCR reaction according to Elshire et al. (2011) original protocol, while for Alfalfa-PV the KAPA library amplification readymix (Kapa Biosystems Cat # KK2611) was used instead. Each library was sequenced in two lanes on Illumina HiSeq 2000 at the Genomic Sequencing and Analysis Facility at the University of Texas at Austin, Texas, USA. For both populations post sequencing analysis and SNP calling was carried out using Tassel UNEAK pipeline (Lu et al. 2013).

Alfalfa is an autotetraploid plant species and can therefore present three different heterozygous genotypes: AAAB, AABB and ABBB. Sequencing depth in this study was not sufficient for accurate tetraploid allelic dosage, but following Li et al. (2014) and Li et al. (2015) reliable genotype calls based on diploid allelic dosage were obtained considering diploid heterozygotes (i.e., AB), while the two homozygous genotypes (AAAA and BBBB) were considered diploid homozygotes (i.e., AA or BB). A further quality filter, implemented through ad hoc Python scripts, removed heterozygous loci with less than 4 and homozygous loci with less than 11 aligned reads. A similar filtering was performed in Rocher et al. (2015) using less restrictive thresholds.

In rice, DNA was isolated from 3-week-old leaves using the DNeasy Plant Mini Kit (QIAGEN) with a TECAN Freedom EVO150 liquid handling robot (TECAN Group Ltd, Switzerland). DNA digestion was performed on 100-ng samples in 96-well plates using ApeKI, which was shown to cut every 1 kb on

Table 1 Descriptive statistics of alfalfa (*M. sativa* L.) and rice (*O. sativa* L.) genotyping

Table 1 Descriptive statistics of alfalfa (<i>M. sativa</i> L.) and rice (<i>O. sativa</i> L.) genotyping		Allowed missing rate per marker (%)	Alfalfa-PV	Alfalfa-Med	Rice
MAF is average minor allele frequency; p(AA), p(AB) and p(BB) are the proportions of AA, AB and BB genotypes. Total number of markers and resulting average missing rate for all markers and for four allowed thresholds of missing rate per markers	100	No. samples	124	154	437
		MAF	0.1724	0.1702	0.140
		p(AA)	0.690	0.691	0.860
		p(AB)	0.274	0.276	–
		p(BB)	0.035	0.032	0.140
		Missing rate	0.666	0.596	0.534
		Markers	32 706	40 734	166 418
		70	Missing rate	0.325	0.29
	Markers		13 190	19 986	109 372
	40	Missing rate	0.161	0.142	0.198
		Markers	7 790	12 931	58 553
	20	Missing rate	0.076	0.074	0.099
		Markers	4 828	8 962	29 872
	10	Missing rate	0.046	0.045	0.049
		Markers	3 405	6 364	15 060

average in a *in silico* digestion of the Nipponbare reference genome, and 96-plex libraries constructed according to the GBS protocol. The libraries were loaded into Genome Analyzer II (Illumina, Inc., San Diego, CA) for sequencing. Raw sequence data filtering, sequence alignment to the rice reference genome (Os-Nipponbare-Reference-IRGSP-1.0, Kawahara et al. 2013), and SNP calling from GBS genotyping were carried out using the Tassel GBS pipeline (Glaubitz et al. 2014).

In total, 32 706, 40 734 and 166 418 SNP markers were called by GBS in Alfalfa-PV, Alfalfa-Med and Rice, respectively (see Table 1). In all datasets SNP variants were renamed so that AA represented the most common homozygote, and BB the least common homozygote.

Imputation methods

We considered six imputation methods: mean value imputation (MNI), K-nearest neighbors imputation (KNNI), Random Forest imputation (RFI), singular value decomposition imputation (SVDI), localized haplotype clustering imputation (“Beagle”) and haplotype reconstruction around the recombination sites (FILLIN from Tassel suite). For all algorithms we imputed a $n \times m$ matrix of n individuals and m markers whose data points, defined in the set $\{0,1,2,NA\}$, represented the three possible genotypes (AA, AB, and BB) and the missing value, respectively.

MNI: in mean value imputation each missing data point was replaced by the mean of the non-missing values for that marker, then discretized to the closest value in $\{0,1,2\}$.

KNNI: in K-nearest neighbors Imputation missing data points were imputed based on the weighted average of the K closest markers (Troyanskaya et al. 2001) defined by the simple matching coefficient distance function (Schwender 2007), specifically designed for categorical data. $K = 4$ was used in KNNI, and neighbors values were weighted by the reciprocal of their distance from the data point to be imputed.

SVDI: Singular value decomposition (SVD) imputation was based on the following factorization of the genotype matrix M (n individuals, m markers):

$$M = U \Sigma V^T \quad (1)$$

where U is a $n \times n$ unitary matrix (i.e., $UU^T = I$), Σ is a $n \times m$ rectangular diagonal matrix of singular values and V^T is the $m \times m$ conjugate transpose of the unitary matrix V . The columns of matrix U and matrix V contain the eigenvectors of $MM^T_{(n \times n)}$ and $M^T M_{(m \times m)}$, respectively, and the corresponding nonzero singular values in Σ are equivalent to the square-root of the non-zero eigenvalues of MM^T and $M^T M$. The first k eigenvectors in U —ordered by decreasing eigenvalue (from Σ)—capture most of the information in the marker genotype matrix M , and were used to generate

linear combinations (principal components) of the original m markers for the imputation of missing data points. The imputation procedure comprised: (1) initial imputation of missing genotypes using MNI, since SVD can only be performed on complete matrices; (2) SVD to select the most informative k eigenvectors of the marker genotype matrix; (3) these k eigenvectors were used as predictors in a linear regression model for marker genotypes:

$$Y = U^* \beta + \varepsilon \quad (2)$$

with Y as vector of n genotypes at marker j ; U^* the $n \times k$ matrix of the first k eigenvectors; β the vector of k regression coefficients; ε the random error terms; (4) all eigenvectors (matrix U^*) and β were used to estimate missing values at marker j . The procedure was repeated [steps (2)–(4)] until convergence. The final imputed data points were then discretized to the closer genotype in $\{0,1,2\}$. A value of $k = 4$ for the eigenvectors to be selected in U^* was used in our implementation of SVDI, based on empirical results in the range 3–20 (data not shown). Additional details on SVDI can be found in Troyanskaya et al. (2001)

RFI: in Random Forest (RF) imputation, missing genotypes at marker j were imputed by means of RF multiple regression trees (Breiman 2001) where all markers other than j were used for the prediction. At each marker j , 100 RF regression trees $\Theta_1 \dots \Theta_{100}$ were grown from a bootstrapped sample of the individuals in Y and a random subset x of $\sqrt{m-1}$ markers. Missing genotypes were imputed averaging predictions over the 100 RF trees:

$$\hat{Y} = \frac{1}{100} \sum_{i=1}^{100} h(x, \Theta_i) \quad (3)$$

where $h(x, \Theta_i)$ is a function of the i_{th} tree and subset of markers. RFI was repeated until convergence or for maximum 10 iterations. After regression, the imputed data were then discretized to the closer genotype in $\{0,1,2\}$.

Beagle: “Localized haplotype clustering imputation” is a method implemented in the software “Beagle” (Browning and Browning 2007). Originally developed for human genetics, Beagle has since found wide application also in animal and plant genetics. Beagle has become so popular that the name of the software is commonly used to refer metonymically to the method it implements, making the two hardly

distinguishable. Beagle infers haplotypes and imputes sporadic missing alleles both with known and unknown phase, using a localized haplotype cluster model. This is a class of directed acyclic graphs which empirically models haplotype frequencies on a local scale and therefore adapts to local structure in the data. Beagle makes use of a hidden Markov model to find the most likely haplotype pair for each individual, given the genotype data for that individual and the graphical haplotype frequency model. The method works iteratively using an expectation-maximization (EM) approach. The imputed missing data, probabilities of missing genotypes and inferred haplotypes are calculated from the model that is fitted in the last iteration. For the imputation of missing genotypes in the completely homozygous rice accessions, a likelihood file (with the prior likelihood of each of the three possible SNP genotypes) was defined to specifically exclude the imputation of non-possible heterozygous genotypes [details in Browning (2011)].

FILLIN The “Fast Inbred Line Library ImputationN” (“FILLIN”) is an imputation method optimized for inbred populations implemented in the “Tassel” programming suite (Swarts et al. 2014). FILLIN is based on haplotype reconstruction around recombination break points. Haplotypes are clustered per genotype similarity using the Hamming distance function. This information is eventually used to impute the target locus in an iterative approach that attempts, through a Markov Chain Monte Carlo (MCMC) process, to maximize the likelihood of the observed SNP calls given the unobserved imputed genotypes. If convergence criteria are not met genotypes are left unimputed. Among the considered algorithms, FILLIN is the only one that can thus produce partially imputed results.

Imputation procedure

In order to assess the imputation performance of the different methods, we introduced increasing proportions of artificial missing genotypes in the data sets. These were then imputed with the various algorithms, measuring the resulting imputation accuracy and computation time. From the overall data, we extracted four datasets that had a maximum of 10, 20, 40 or 70 % missing data per marker. For each of these datasets, we randomly introduced an additional 1, 5, 10 or 20 % missing genotypes, on which imputation accuracy could be measured. Table 1 reports the

number of markers and average missing rate (before introducing artificial missing genotypes) for the four call rate thresholds. This procedure produced 16 combinations for each population: 4 call rates \times 4 levels of artificial missing markers. On each combination, we applied the 6 described imputation methods. Rice data were analyzed as a whole and by each of the 12 chromosomes separately. In absence of the alfalfa genome, the *M. Truncatula* genome was used as reference for those algorithms (Beagle and FILLIN) requiring marker position.

To further investigate the effect of presence (or absence) of genome information we generated “reshuffled” datasets where marker positions were randomly assigned inside each chromosome. This way, the linkage disequilibrium between markers is broken and the relative performance decrease can be directly assessed. The reshuffled datasets were tested only on Beagle and FILLIN.

Imputation accuracy and efficiency assessment

We used the artificial missing genotypes to measure the performance of the six imputation methods. The rice dataset contained only two genotype classes: AA (homozygous for the major allele) and BB (homozygous for the minor allele). Both alfalfa datasets contained also the third genotype classes AB (heterozygous, pooling the three possible heterozygotes AAAB, AABB, AB BB).

For each experiment we measured the overall imputation accuracy and the imputation accuracy within each genotype class. The imputation accuracy was computed as the number of missing data correctly imputed divided by the total number of artificially missing data (proportion of correct imputations):

$$\text{accuracy} = \frac{1}{n} \sum_{i=1}^n I(g_i = \hat{g}_i) \quad (4)$$

where $I()$ is an indicator function that returns 1 if the imputed (\hat{g}) and true (g) genotypes are equal, 0 otherwise. We obtained four accuracy measures for alfalfa and three for rice.

For each imputation method, the total computation time needed to complete the analyses was measured as an indicator of their relative performance. The measured time equals the total dedicated CPU time in case of single thread execution, and corresponds to a

fraction of it in case of algorithm supporting multi-thread execution. In our experiments only RFI implemented parallel execution, and it was configured to use 10 CPUs: thus, all RFI reported times are to be considered as time elapsed while employing ten times more computational resources compared to the other algorithms.

To ensure consistency in the experimental conditions, all analysis were run on the same computing platform at PTP Science Park (www.ptp.it): a multi-node cluster with 64 CPUs AMD Opteron(tm) 2400 MHz and 256 GB RAM for each node.

Software

Data handling, editing and preparation, and summary of results and plots were all performed using the open-source environment for statistical programming R Core Team (2014). All imputation methods except Beagle and FILLIN were implemented in R: MNI directly in base R; KNNI using the function *knncatimpute()* from the ScRime package (Schwender and Fritsch 2013); SVDI with the R package “bcv” (Perry 2009); RFI using the “MissForest” (Stekhoven and Bhlmann 2012) R package (with parameters: *ntree*=100, *maxiter*=10, *parallelize*=‘variables’). The “Beagle” localized haplotype clustering imputation method was implemented with the Beagle software version 3.3.2 (Browning and Browning 2007). The FILLIN algorithm (Swarts et al. 2014) is implemented using the Tassel CLI plugin *FILLINindHaplotypesPlugin* followed by *FILLINImputationPlugin*. The latter allowed accuracy measurements through *-accuracy* and *-propSitesMask* options.

The Bowtie 2 tool (Langmead and Salzberg 2012) was used to query the consensus sequence of each alfalfa tag pair containing a SNP against the *M. truncatula* reference genome Version 4.0 using the *-verySensitivelocal* preset. The BWA alignment tool (Li and Durbin 2009) was used to align rice tags on Rice Genome Annotation Project Release 7.

Results

Genotypes

GBS detected 32 706, 40 734 and 166 418 SNP loci in Alfalfa-PV, Alfalfa-Med and Rice, respectively (Table 1). The overall missing rate was 66.6, 59.6 and 53.4 % in the

three datasets. Supplementary Figure S1 shows the density distribution of missing rate per marker in the complete datasets. The amount of missing marker genotypes varied with the threshold of maximum per-marker missing rate allowed in the data (10, 20, 40 and 70 %) and the proportion of artificial missing genotypes that were introduced (1, 5, 10 and 20 %). Missing data point counts ranged from a minimum of 3 878 in Alfalfa-PV with 10 % allowed and 1 % artificial missing genotypes, to a maximum of 5 580 616 in Rice with 70 % allowed and 20 % artificial missing genotypes. Over all datasets and allowed/artificial missing genotypes thresholds, the amount of missing data points to be imputed averaged 249 405.

Minor Allele Frequency (MAF) was 0.172, 0.170 and 0.140 in Alfalfa-PV, Alfalfa-Med and Rice, respectively. The three datasets were therefore unbalanced with respect to genotype classes: Rice had the lowest MAF, while Alfalfa-PV and Alfalfa-Med had the smallest minor classes (minor homozygotes 3.5 and 3.2 % of the total genotypes). Figure 1 reports the proportion of genotypes in each class in the complete datasets (100 %) and as a function of the maximum allowed missing rate per marker. The relative proportion of genotype classes remained approximately stable in rice. In contrast, the heterozygous class in alfalfa tended to become smaller with increasing proportion of allowed genotypes. This reflected the different GBS calling criteria for homozygous and heterozygous SNP: heterozygous SNPs required fewer overlapping reads to be called, while the implemented quality filters on alfalfa required a larger number of reads to call a homozygous locus. Therefore homozygous loci tended to have higher missing rates, and to be included progressively more frequently with more relaxed thresholds on allowed missing level. The average missing rate and MAF in the 12 rice chromosomes were comparable to the whole-rice dataset.

Alignment on reference genomes assigned a position to 88.66, 57.54 and 57.86 % of Rice, Alfalfa-Med and Alfalfa-PV markers, respectively, reflecting the difference between having a reference genome available (rice) or using that of a closely related species (alfalfa).

Accuracies are reported, averaged, in Supplementary Table 1. In rice, the average imputation accuracy over all genotype classes was above 80 % for all methods. No significant difference was found between chromosomes. Thus average accuracy over all chromosomes is therefore presented (Fig. 2).

The overall accuracy for alfalfa, averaged across the two datasets, was 12–25% lower than rice, ranging between 66 and 82 %. No significant differences were found between the two datasets. Thus average accuracy over the Alfalfa-Med and Alfalfa-PV is presented (Fig. 3).

The imputation accuracy varied across the different genotype classes, with the highest accuracy in the most common genotype (averaging 92.27 and 96.47 % in alfalfa and rice, respectively), while the least common genotype class showed much lower accuracy (averaging 5.79 and 69.88 % in alfalfa and rice, respectively). Alfalfa datasets included heterozygous SNP loci, accounting for an average of 27 % of individuals per locus (Table 1). Heterozygotes had an intermediate imputation accuracy ranging from 0.5 % with Beagle to 66.03 % with KNNI. Results for each individual alfalfa population and the 12 individual rice chromosomes are similar to those discussed here (see Supplementary Figures S2 to S15).

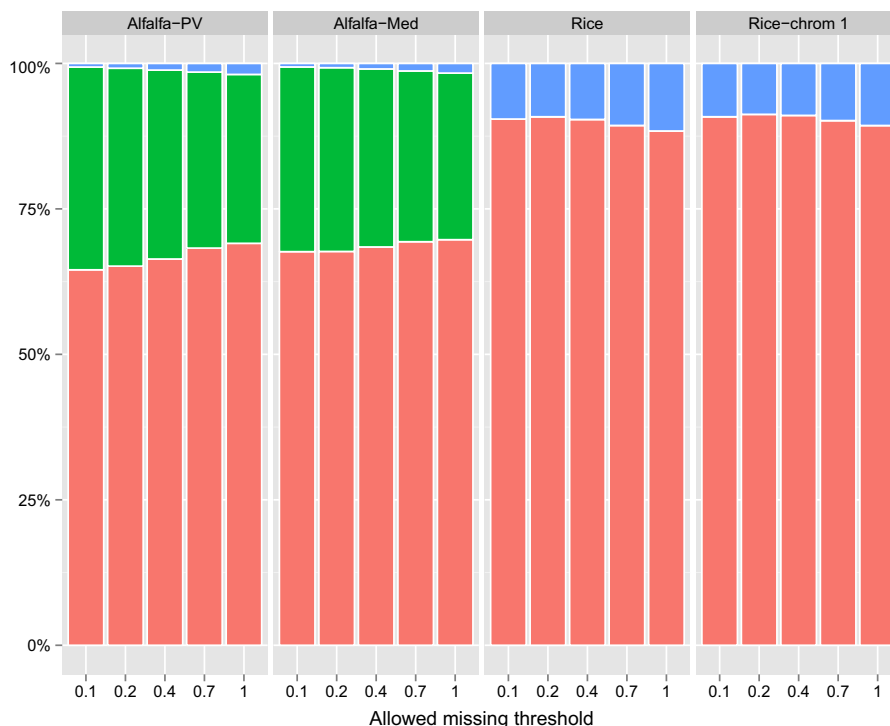
Missing rate did not affect the imputation accuracy substantially, with most algorithms showing a flat or almost flat response to increased missing rate. KNNI showed a decreasing imputation accuracy with increasing missing rate, both in alfalfa and rice.

Beagle outperformed all other imputation methods in rice, making efficient use of marker position (Fig. 2). When markers were randomly reshuffled (thus offsetting position information), though, general imputation methods (other than MNI) showed a higher imputation accuracy than Beagle. The accuracy loss with shuffled markers is exacerbated in alfalfa datasets, and in particular in heterozygous and minor homozygous genotypes.

Similarly, FILLIN resulted in a 5–8 % accuracy drop when marker positions were shuffled. FILLIN performances were in general lower than other methods due to the fraction of unimputed genotypes, here accounted as imputation errors. On average, FILLIN did not impute 14.6 % of missing data in Rice and practically the entirety of Alfalfa data. Even trying to relax FILLIN parameters, we couldn't obtain meaningful results in alfalfa. Thus, FILLIN accuracy results are not reported in Fig. 3.

When examining general imputation methods, resulting accuracies highlighted RFI and KNNI as the best ones. The two methods had comparable performances over all experiments. SVDI followed, and MNI ranked last.

Fig. 1 Proportion of genotype classes for the different missing rate thresholds (0.1, 0.2, 0.4, 0.7 and 1 -the complete dataset), in Alfalfa-PV, Alfalfa-Med and Rice (all chromosomes and chromosome 1 only). Red bars represents the most common homozygous (AA), blue bars the least common homozygous (BB) and green bars the heterozygous (AB). (Color figure online)



In Rice, for 70 % allowed missing rate, the imputation accuracy of RFI was still higher than 90 %, while it dropped toward 85 % for KNNI and SVDI. As expected, MNI had an imputation accuracy of 100 % in the major homozygous class, but dropped to 0 % in the minor homozygous class.

In alfalfa KNNI showed the highest imputation accuracy over all classes (averaging 82.20 %) and in the heterozygous (66.03 %) and minor homozygous (14.68 %) genotype classes (Fig. 3). RFI and SVDI were second ranking. They did not differ much from MNI in the overall accuracy (MNI: 78.38 %, RFI: 79.50 %, SVDI: 77.36 %) and in the major (MNI: 88.69 %, RFI: 89.05 %, SVDI 85.81 %) homozygous class, but gave higher imputation accuracy in the heterozygous (MNI: 58.88 %, RFI: 61.75 %, SVDI 61.67 %). In the minor homozygous class only SVDI resulted in some sporadic correct imputation, while RFI and MNI did not produce any sensible result (MNI: 0 %, RFI: 0.01 %, SVDI 4.47 %) classes.

Computation time

The amount of time required to complete the imputation process was recorded for each method. Only the implementation of RFI could leverage a multi-

core/multi-thread environment, so that RFI computation times should be evaluated considering 10 CPUs used in parallel, while all other algorithms used only one CPU at a time.

Imputation efficiency (Fig. 4) was assessed with respect to the gross dimension of the data set (i.e., number of markers \times number of samples). An alternative analysis relating computation times to the number of missing genotypes brought to the same results (data not shown).

RFI required by far the longest computation times (in spite of parallelization), which grew approximately exponentially ($O(e^N)$) with problem size, and easily required tens of hours for individual rice chromosomes, and hundred of hours (up to 937 h) for the complete rice dataset. The second slowest algorithm was KNNI, with computation times growing approximately quadratically ($O(N^2)$) with problem complexity (the chosen KNNI implementation contains no heuristic method to prune the all vs. all distance calculation). KNNI was however significantly faster than RFI (on average 16.7 times faster), requiring a time ranging from 20.68 s (rice chromosome 7, 10 % allowed, 1 % artificial) to about 28 h (rice complete dataset, 70 % allowed, 20 % artificial) to complete the imputation task.

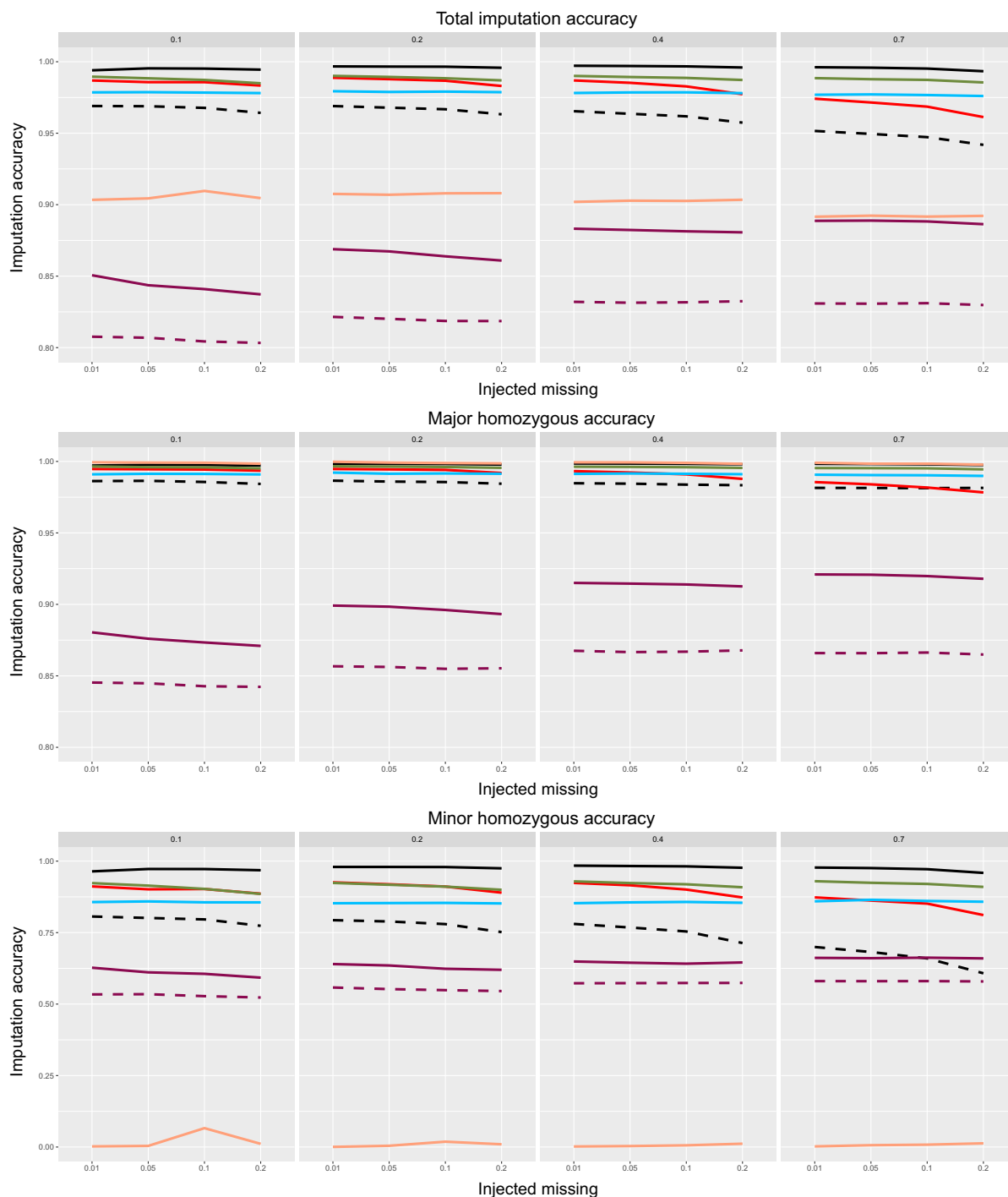


Fig. 2 Imputation accuracies overall, for the major homozygous genotype (AA), and for the minor homozygous genotype (BB) in datasets consisting of 10, 20, 40 and 70 % allowed missing data per locus (*boxes*) with 1, 5, 10 and 20 % additional missing values artificially introduced (*x-axis*) averaged over the 12 rice chromosomes. *Lines colors* represent the five imputation

algorithms: MNI (*salmon*), KNNI (*red*), SVDI (*blue*), RFI (*green*), Beagle with ordered markers (*solid black*), Beagle with unordered markers (*dashed black*), FILLIN with ordered markers (*purple*) and FILLIN with unordered markers (*dashed purple*). *Y axis scale* changes to highlight differences. (Color figure online)

All other algorithms were faster, with computation times growing linearly or logarithmically with problem size. Beagle, FILLIN and SVDI resulted in similar execution times. MNI was by far the fastest imputation algorithm, being scarcely affected by the size of the problem (computation times ranging from 0.23 to 31.29 s).

Discussion

Minor allele frequency and data (un)balancedness

GBS data pose a greater challenge than SNP array data to imputation algorithms, mainly as a consequence of the much larger quantity of missing genotypes they contain. In our data sets we found that missing rates varied from 53 to 67 %. With larger amounts of missing data the complexity of the imputation problem increases. Imputation errors can have a negative impact on successive analyses [e.g., genomic predictions Rutkoski et al. (2013); Annicchiarico et al. (2015)].

The imputation of missing SNP genotypes is a special case of the broader family of classification problems: three-class missing genotypes (AA, AB and BB) at any given SNP locus are classified based on known genotypes at all remaining data points. Classification problems are known to be harder when data are unbalanced, i.e., the classes appear at different frequencies in the datasets, with typically one over-represented class (see Kotsiantis et al. 2006; Sun et al. 2009 for a review). In SNP genotype imputation, data balancedness is directly related to the minor allele frequency (MAF). In the classification of unbalanced observations, it is important to look not only at the total classification accuracy, but also at the per-class accuracy. The total classification accuracy may be misleading, being “dominated” by the majority class (He and Garcia 2009). Indeed, we found that for most methods, even when the total classification accuracy was very high, relatively large error rates were present in the minority classes. The dependency of imputation results on MAF has already been acknowledged (e.g., Hickey et al. (2012) in maize; Ma et al. (2013) in cattle; Pei et al. (2008) in humans). Our per-class dissection of results allowed a deeper insight into the imputation process. Indeed, all imputation algorithms performed considerably better in the majority class

Fig. 3 Imputation accuracies overall, for the major homozygous genotype (AA), for heterozygotes (AB), and for the minor homozygous genotype (BB) in datasets consisting of 10, 20, 40 and 70 % allowed missing data per locus (*boxes*) with 1, 5, 10 and 20 % additional missing values artificially introduced (*x*-axis) averaged two alfalfa populations (Alfalfa-Med and Alfalfa-PV). *Lines colors* represent the five imputation algorithms: MNI (*salmon*), KNNI (*red*), SVDI (*blue*), RFI (*green*), Beagle with ordered markers (*solid black*) and Beagle with unordered markers (*dashed black*). FILLIN was unable to impute alfalfa data and is absent from figure. *Y* axis scale changes to highlight differences. (Color figure online)

rather than in the less frequent classes. In alfalfa, KNNI was easily the best imputation method in the heterozygous and minor homozygous classes, while SVDI and RFI performed only slightly better than MNI. In rice, Beagle gave the best imputation results, with accuracy close to 100 % both in the major and minor homozygous classes.

Missing rate and the “curse of dimensionality”

In general, imputation methods were relatively robust to increasing missing rates. Only KNNI, and to a lesser extent RFI, showed slowly degrading performances at very high missing rates. The divergent response of the different imputation methods to missing rate became apparent only in the most challenging scenarios, when markers with up to 40–70 % missing genotypes were allowed in the dataset.

KNNI’s decrease in imputation accuracy with increasing missing rates can be interpreted in relation to the phenomenon known as the “curse of dimensionality” (Bellman 1957; Marimont and Shapiro 1979): with the same amount of data, the increasing number of parameters—the “dimensions” of the problem - increases and complicates the identification of k neighbors which are close enough to the data point to be classified/imputed. On average, the side l of the hypercube to include k neighbors is a function of k , n (sample size) and p (number of parameters): $l = \left(\frac{k}{n}\right)^{1/p}$. With n constant, the hypercube in which the k neighbors lie gets bigger as p increases. This holds especially for local learning methods (e.g., K-nearest neighbors methods, local regression) which rely heavily on the information content of the neighborhood. RFI—which only partially relies on local structure of the data—suffered marginally from high

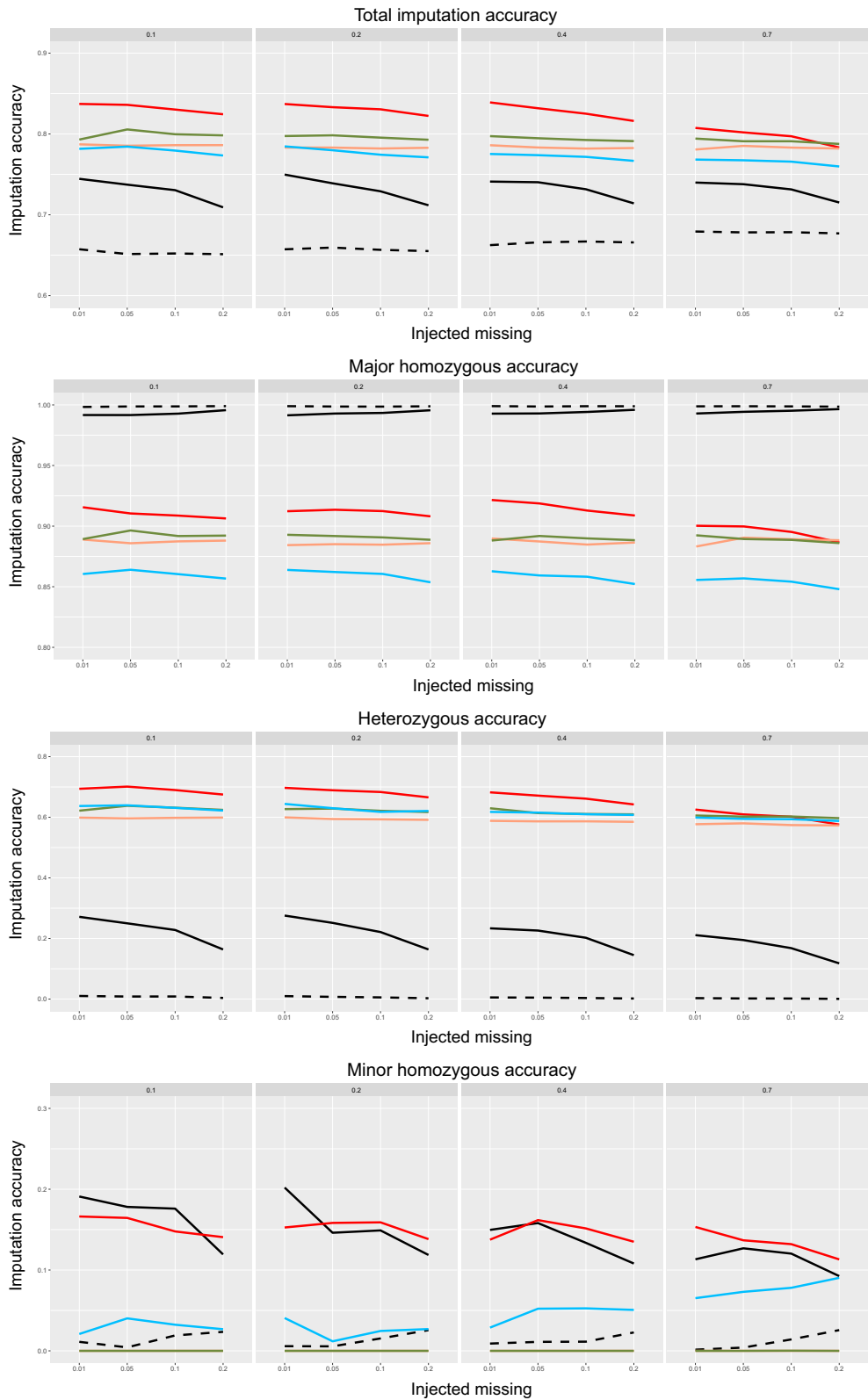
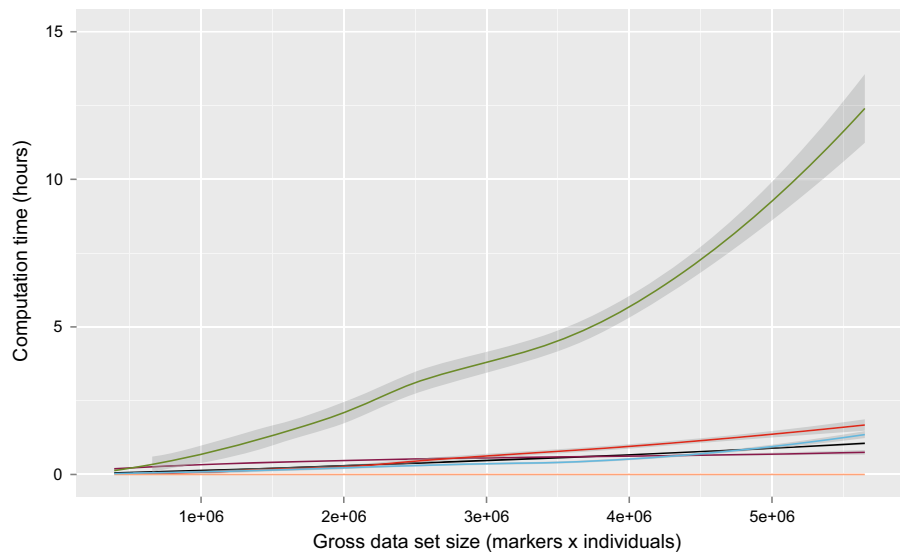


Fig. 4 computation time as a function of the total size of the imputed dataset. *Lines colors* represent the five imputation algorithms: MNI (*salmon*), KNNI (*red*), SVDI (*light blue*), RFI (*green*), Beagle (*black*) and FILLIN (*purple*). Plots include measures on Alfalfa-Med population, Alfalfa-PV population, and rice chromosomes 1–12. Complete rice datasets are omitted for readability. (Color figure online)



missingness. SVD and MNI, which build on more general properties of the data, and Beagle—whose learning process is fundamentally different and specific for genotype imputation—were scarcely affected by missing rates.

Imputation efficiency: differences in rice versus alfalfa

Imputation results were significantly different in rice and alfalfa: all imputation algorithms performed consistently better in rice, where Beagle, RFI and KNNI achieved performances comparable to what is reported in literature for SNP chip data (e.g., VanRaden et al. (2013) for bovine data, The 1000 Genomes Project Consortium (2012) for human data). On the other hand, imputation accuracy in alfalfa was much lower. This can be ascribed to four main factors:

- (1) The imputation problem was simpler in rice than in alfalfa, since the rice dataset comprised only two genotype classes (AA and BB) instead of alfalfa's three (AA, AB, BB);
- (2) Rice is natively diploid while alfalfa, autotetraploid, has been rendered diploid “in silico” during SNP calling. This simplification step made alfalfa data less adherent to the underlying biology;
- (3) Rice data have higher marker density (2.4 Kbp/SNP, compared to 24.5 Kbp/SNP for Alfalfa-

PV and 19.6 Kbp/SNP for Alfalfa-Med) due to both a higher number of markers and a smaller genome (400 Mb for rice, 800 Mb for Alfalfa). This allowed for higher average values of linkage disequilibrium (LD) among SNP markers, thus facilitating the imputation process;

- (4) Rice markers were aligned on their native genome, while alfalfa markers were aligned on the genome of a different species.

There was no population structure in alfalfa (Annicchiarico et al. 2015), while the rice dataset was intrinsically stratified—being a collection of five subpopulations (*indica*, *temperate japonica*, *tropical japonica*, *aus* and *aromatic*). Most imputation methods implicitly exploit population structure (e.g., KNNI computes distances based on genetic relatedness; Beagle reconstructs haplotypes based on genetic similarities), without explicitly modelling it. The imputation of missing genotypes, however, is not an inferential problem such as GWAS (genome-wide association study), where the significance and estimate of SNP effects are known to be inflated if population structure is not included in the model, and the impact on the accuracy of imputation is likely to be small. In cattle from Scandinavian countries, (Brøndum et al. 2012) actually found higher imputation accuracy when combining populations (cattle breeds) in Beagle (without explicitly modelling the population structure), most likely as a result of the larger sample size.

On the other hand, population stratification may help explain why FILLIN performed poorly in our rice dataset, which is a collection of subpopulations, while this algorithm is optimized for homogeneous inbred populations. However, the influence of population structure on the accuracy of imputation was not formally tested in this work, and this remains an interesting topic for further investigation.

Reference genome assembly and ordered versus unordered markers

When a reference genome is available, marker position can be used in the imputation process. All methods specifically developed for the imputation of missing genotypes have been designed for markers ordered along the genome sequence, and make explicit or implicit use of position information.

In rice, where a reference genome is available, we could assess the effect of using ordered vs. unordered markers for imputation. For alfalfa there is no reference genome sequenced yet. The genome of the close relative diploid *M. Truncatula* is available (Young et al. 2011) and can in principle be used. However, while the two genomes show high synteny (Li et al. 2014), aligning on a different species (and with different ploidity) comes at the price of discarding those markers that do not align. In our case only 57.54 % (Alfalfa-Med) and 57.86 % (Alfalfa-PV) aligned on the *M. Truncatula* genome, compared to the 88.66 % of rice markers aligning on *O. Sativa* genome. This left with 23 438 and 18 923 SNPs to be used for analysis with ordered markers.

Among the imputation methods used in this study, Beagle and FILLIN are the only ones that make use of marker information, and were therefore tested with ordered and randomly shuffled markers. In rice dataset with ordered markers Beagle showed astounding resilience to high missingness, with imputation accuracy in the minority class >95 % even in the most extreme scenarios (70 % allowed, 20 % artificial missing genotypes). Beagle worked substantially worse when the marker position was randomly reshuffled, with imputation accuracy in the minority class ranging from 75 % (10 % allowed, 1 % artificial missing rates) to almost 50 % (70 % allowed, 20 % artificial missing rates). When tested with alfalfa Beagle largely overestimated the majority class and had very poor performances on the heterozygous class.

Shuffling marker positions further worsened Beagle performance in alfalfa, and almost all missing genotypes were assigned to the majority class.

FILLIN performance in rice followed the same pattern and was substantially worse when markers were shuffled. This result confirmed that imputation methods specifically designed for ordered markers are not an option for species lacking a reference genome.

Previous studies, though not detailing a per-class breakdown of the imputation accuracy, provide interesting comparisons. In sugar beet, Biscarini et al. (2014) used Beagle to impute SNP chip markers with partial within-chromosome/scaffold alignment, finding global accuracies ranging from 84 to 80.9 % with 1 to 20 % missing genotypes. Huang et al. (2014) tested several algorithms, including Beagle, on simulated ordered GBS rice data with missing rates in the 30–60 % range, and found that Beagle gave imputation accuracy consistently higher (+15–20%) than KNNI. Finally, Swarts et al. (2014) used ordered GBS maize data to compare Beagle with the FSFHap and FILLIN algorithms they developed. They mostly obtained high imputation accuracies in the three genotype classes, providing further evidence of the added value of having a reference genome assembly.

Size of the imputation problem

Albeit this was not the main objective of the paper, and no formal effort to optimize the implementation of the various imputation methods was done, still the recorded computation times give us some interesting information. The size of the problem, besides affecting the accuracy achievable by some imputation methods, is mainly relevant with reference to the required computation resources. RFI was the most demanding algorithm, with computation times growing exponentially both with the total number of marker genotypes ($m \text{ SNP} \times n \text{ samples}$) and with the number of missing genotypes to be imputed. When analyzing the complete rice dataset (all chromosomes together), RFI took a maximum of about 40 days to complete imputation, and became computationally intractable (on the available platform) for the largest missing rate scenarios. This was partially mitigated by parallelization (we used 10 CPUs in our experiments). All other imputation algorithms took much shorter times to complete imputation, even in the most challenging scenarios. Putting together imputation

accuracy and computation time, we found that the best performing imputation algorithms were Beagle in rice and KNNI in alfalfa.

Conclusions

Imputation of missing genotypes can be an effective technique also for GBS data. The most accurate imputation methods resulted in a total amount of wrongly imputed genotypes near to zero in rice and slightly over 10 % in alfalfa. The proportion of imputation errors, however, varied dramatically among genotype classes, approaching very high levels with the worst methods on the minor homozygous class. The structure of the genome and the maturity of the reference assembly play a role in the imputation efficiency, as indicated by the greater efficiency obtained in rice compared with alfalfa. While Beagle was preferable for rice, general data imputation methods performed significantly better in the absence of a reference genome. In particular, in alfalfa RFI and KNNI showed the lowest imputation errors in all classes, with KNNI to be preferred for computational efficiency. The alignment of markers to closely related species can help, but comes at the price of discarding markers that do not align, and with general methods still performing better.

The imputation accuracy tended to be relatively robust over increasing missing rates. However, KNNI showed a somewhat lower accuracy on high missingness scenarios, probably as a consequence of the “curse of dimensionality”. In terms of computation requirements, all methods proved to be tractable over all problem complexities with a standard bioinformatics-lab computation infrastructure, except RFI (whose computational requirements increased exponentially with problem size and took up to 40 days with the complete rice dataset and maximum missing rate thresholds).

The results of this paper showed that high imputation accuracies can be achieved with GBS data, and that general imputation methods are a valid option when the reference genome is not available or the alignment on the genome of a closely related species leads to losing too many markers. Additionally, we highlighted the importance of examining imputation accuracy in the different genotype classes, since the common practice of summarizing performances in a

single index can be very misleading when data are unbalanced (i.e., MAF is low). Exploring ways to improve imputation accuracy in GBS data is an important research area and can help making genotyping-by-sequencing a very attractive and cost-effective technique for genomic applications in species with and without a reference genome.

Acknowledgments The rice data used in this research paper were produced within the framework of the Italian national project “RISINNOVA” (Grant No. 2010–2369), financially supported by the AGER Foundation. The creation of the alfalfa data sets was funded by the projects Genomic selection in alfalfa (GENALFA) funded by the Italian Ministry of Foreign Affairs and International Cooperation in the framework of the Italy-USA scientific cooperation program, the Italian share of the FP7-ArimNet project Resilient, water- and energy-efficient forage and feed crops for Mediterranean agricultural systems (REFORMA) funded by the Italian Ministry of Agricultural and Forestry Policies.

References

- Annicchiarico P, Nazzicari N, Li X, Wei Y, Pecetti L, Brummer EC (2015) Accuracy of genomic selection for alfalfa biomass yield in different reference populations. *BMC Genomics* 16(1):1–13. doi:[10.1186/s12864-015-2212-y](https://doi.org/10.1186/s12864-015-2212-y)
- Aulchenko YS, Ripke S, Isaacs A, Van Duijn CM (2007) GenABEL: an R library for genome-wide association analysis. *Bioinformatics* 23(10):1294–1296
- Bellman R (1957) *Dynamic programming*. Princeton University Press, Princeton
- Biscarini F, Stevanato P, Broccanello C, Stella A, Saccomani M (2014) Genome-enabled predictions for binomial traits in sugar beet populations. *BMC Genet* 15(1), 87. <http://www.biomedcentral.com/1471-2156/15/87/>
- Breiman L (2001) Random forests. *Mach Learn* 45(1):5–32. <http://link.springer.com/article/10.1023/A:1010933404324>
- Brøndum RF, Ma P, Lund MS, Su G (2012) Short communication: Genotype imputation within and across nordic cattle breeds. *J Dairy Sci* 95(11):6795–6800
- Browning SR, Browning BL (2007) Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *The Am J Hum Genet* 81(5):1084–1097. doi: [10.1086/521987](https://doi.org/10.1086/521987). <http://www.sciencedirect.com/science/article/pii/S0002929707638828>
- Browning B (2011) Beagle 3.3.2. https://faculty.washington.edu/browning/beagle/beagle_3.3.2_31Oct11.pdf
- Crossa J, Beyene Y, Kassa S, Prez P, Hickey JM, Chen C, Campos Gdl, Burgueo J, Windhausen VS, Buckler E, Jannink JL, Cruz MAL, Babu R (2013) Genomic prediction in maize breeding populations with genotyping-by-sequencing. *G3: Genes|Genomes|Genetics* 3:11:1903–1926. doi: [10.1534/g3.113.008227](https://doi.org/10.1534/g3.113.008227). <http://www.g3journal.org/content/3/11/1903>

- Elshire RJ, Glaubitz JC, Sun Q, Poland JA, Kawamoto K, Buckler ES, Mitchell SE (2011) A robust, simple genotyping-by-sequencing (GBS) approach for high diversity species. *PLoS One* 6(5):e19379. doi:[10.1371/journal.pone.0019379](https://doi.org/10.1371/journal.pone.0019379)
- Endelman JB (2011) Ridge regression and other kernels for genomic selection with r package rrblup. *Plant Genome* 4:250–255
- Glaubitz JC, Casstevens TM, Lu F, Harriman J, Elshire RJ, Sun Q, Buckler ES (2014) TASSEL-GBS: a high capacity genotyping by sequencing analysis pipeline. *PLoS One* 9(2):E90346. <http://dx.plos.org/10.1371/journal.pone.0090346>
- Hayes B, Bowman P, Chamberlain A, Goddard M (2009) Invited review: genomic selection in dairy cattle: progress and challenges. *J Dairy Sci* 92(2):433–443
- He H, Garcia EA (2009) Learning from imbalanced data. *IEEE Trans Knowl Data Eng* 21(9):1263–1284
- Hickey JM, Crossa J, Babu R, de los Campos G (2012) Factors affecting the accuracy of genotype imputation in populations from several maize breeding programs. *Crop Sci* 52:2:654 doi:[10.2135/cropsci2011.07.0358](https://doi.org/10.2135/cropsci2011.07.0358). <https://www.crops.org/publications/cs/abstracts/52/2/654>
- Huang BE, Raghavan C, Mauleon R, Broman KW, Leung H (2014) Efficient imputation of missing markers in low-coverage genotyping-by-sequencing data from multi-parental crosses. *Genetics* 197(1):401–404. doi:[10.1534/genetics.113.158014](https://doi.org/10.1534/genetics.113.158014). <http://www.genetics.org/content/197/1/401>
- International Rice Genome Sequencing Project (2005) The map-based sequence of the rice genome. *Nature* 436:7052:793–800. <http://www.nature.com/articles/nature03895>
- Kawahara Y, de la Bastide M, Hamilton JP, Kanamori H, McCombie WR, Ouyang S, Schwartz DC, Tanaka T, Wu J, Zhou S, et al (2013) Improvement of the *Oryza sativa* Nipponbare reference genome using next generation sequencing and optical map data. *Rice* 6(1):4. <http://www.biomedcentral.com/content/pdf/1939-8433-6-4.pdf>
- Kotsiantis S, Kanellopoulos D, Pintelas P (2006) Handling imbalanced datasets: a review. *GESTS Int Trans Comput Sci Eng* 30(1):25–36
- Langmead B, Salzberg SL (2012) Fast gapped-read alignment with Bowtie 2. *Nature Methods* 9(4):357–359. <http://www.nature.com/nmeth/journal/v9/n4/abs/nmeth.1923.html>
- Li H, Durbin R (2009) Fast and accurate short read alignment with burrows-wheeler transform. *Bioinformatics* 25(14):1754–1760
- Li X, Wei Y, Acharya A, Hansen JL, Crawford JL, Viands DR, Michaud R, Claessens A, Brummer EC (2015) Genomic prediction of biomass yield in two selection cycles of a tetraploid alfalfa breeding population. *Plant Genome*. doi:[10.3835/plantgenome2014.12.0090](https://doi.org/10.3835/plantgenome2014.12.0090). <https://www.crops.org/files/publications/tpg/first-look/plantgenome-tpg-2014-12-0090.pdf>
- Li X, Wei Y, Acharya A, Jiang Q, Kang J, Brummer EC (2014) A saturated genetic linkage map of autotetraploid alfalfa (*Medicago sativa* L.) developed using genotyping-by-sequencing is highly syntenous with the *Medicago truncatula* genome. *G3: Genes Genomes Genetics* 4(10):1971–1979 (2014). <http://www.g3journal.org/content/4/10/1971.short>
- Lu F, Lipka AE, Glaubitz J, Elshire R, Cherney, JH, Casler MD, Buckler ES, Costich DE Switchgrass genomic diversity, ploidy, and evolution: novel insights from a network-based snp discovery protocol. *PLoS Genet* 9(1):e1003215. doi:[10.1371/journal.pgen.1003215](https://doi.org/10.1371/journal.pgen.1003215)
- Ma P, Brndum RF, Zhang Q, Lund MS, Su G (2013) Comparison of different methods for imputing genome-wide marker genotypes in Swedish and Finnish Red Cattle. *J Dairy Sci* 96(7):4666–4677. <http://www.sciencedirect.com/science/article/pii/S0022030213003664>
- Marimont RB, Shapiro MB (1979) Nearest neighbour searches and the curse of dimensionality. *IMA J Appl Math* 24(1):59–70. doi:[10.1093/imamat/24.1.59](https://doi.org/10.1093/imamat/24.1.59). <http://imamat.oxfordjournals.org/content/24/1/59>
- Nicolazzi EL, Biffani S, Biscarini F, Orozco ter Wengel P, Caprera A, Nazzicari N, Stella A (2015) Software solutions for the livestock genomics SNP array revolution. *Anim Genet*. doi:[10.1111/age.12295](https://doi.org/10.1111/age.12295). <http://onlinelibrary.wiley.com/doi/10.1111/age.12295/abstract>
- Pei YF, Li J, Zhang L, Papasian CJ, Deng HW (2008) Analyses and comparison of accuracy of different genotype imputation methods. *PLoS One* 3(10):e3551. <http://dx.plos.org/10.1371/journal.pone.0003551>
- Pérez P, de los Campos G (2014) Genome-wide regression & prediction with the bgrr statistical package. *Genetics pp. genetics*–114
- Perry PO (2009) Bcv: cross-validation for the SVD (bi-cross-validation). <http://cran.r-project.org/web/packages/bcv/index.html>
- Poland J, Endelman J, Dawson J, Rutkoski J, Wu S, Manes Y, Dreisigacker S, Crossa J, Sanchez-Villeda H, Sorrells M, Jannink JL (2012) Genomic selection in wheat breeding using genotyping-by-sequencing. *Plant Genome* 5(3):103. doi:[10.3835/plantgenome2012.06.0006](https://doi.org/10.3835/plantgenome2012.06.0006). <https://www.crops.org/publications/tpg/abstracts/5/3/103>
- R Core Team: R (2014) A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. <http://www.R-project.org>
- Rocher S, Jean M, Castonguay Y, Belzile F (2015) Validation of genotyping-by-sequencing analysis in populations of tetraploid alfalfa by 454 sequencing. *PLoS One* 10(6):e0131918. doi:[10.1371/journal.pone.0131918](https://doi.org/10.1371/journal.pone.0131918)
- Rutkoski JE, Poland J, Jannink JL, Sorrells ME (2013) Imputation of unordered markers and the impact on genomic selection accuracy. *G3: Genes Genomes Genetics* 3(3):427–439. <http://www.g3journal.org/content/3/3/427.short>
- Schwender H (2007) Statistical analysis of genotype and gene expression data. Ph.D. thesis. <https://eldorado.tu-dortmund.de/handle/2003/23306>
- Schwender H, Fritsch A (2013) Scribe: analysis of high-dimensional categorical data such as SNP data. <http://cran.r-project.org/web/packages/scribe/index.html>
- Stekhoven DJ, Bhlmann P (2012) MissForest non-parametric missing value imputation for mixed-type data. *Bioinformatics* 28(1):112–118. <http://bioinformatics.oxfordjournals.org/content/28/1/112.short>
- Sun Y, Wong AK, Kamel MS (2009) Classification of imbalanced data: a review. *Int J Pattern Recognit Artif Intell* 23(04):687–719. doi:[10.1142/S0218001409007326](https://doi.org/10.1142/S0218001409007326). <http://www.worldscientific.com/doi/abs/>
- Swarts K, Li H, Romero Navarro JA, An D, Romay MC, Hearne S, Acharya C, Glaubitz JC, Mitchell S, Elshire RJ, Buckler ES, Bradbury PJ (2014) Novel Methods to optimize

- genotypic imputation for low-coverage, next-generation sequence data in crop plants. *Plant Genome* 7(3):0. doi:10.3835/plantgenome2014.05.0023. <https://www.crops.org/publications/tpg/abstracts/7/3/plantgenome2014.05.0023>
- The 1000 Genomes Project Consortium (2012) An integrated map of genetic variation from 1,092 human genomes. *Nature* 491(7422):56–65. doi:10.1038/nature11632
- Troyanskaya O, Cantor M, Sherlock G, Brown P, Hastie T, Tibshirani R, Botstein D, Altman RB (2001) Missing value estimation methods for DNA microarrays. *Bioinformatics* 17(6):520–525 (2001). <http://bioinformatics.oxfordjournals.org/content/17/6/520.short>
- VanRaden PM, Null DJ, Sargolzaei M, Wiggans GR, Tooker ME, Cole JB, Sonstegard TS, Connor EE, Winters M, vanKaam JBCHM, Valentini A, Van Doormaal BJ, Faust MA, Doak GA (2013) Genomic imputation and evaluation using high-density Holstein genotypes. *J Dairy Sci* 96(1):668–678 (2013). doi:10.3168/jds.2012-5702. <http://www.science.direct.com/science/article/pii/S0022030212007576>
- VanRaden PM, O'Connell JR, Wiggans GR, Weigel KA (2011) Genomic evaluations with many more genotypes. *Genet Sel Evol* 43(10):10–1186. <http://www.biomedcentral.com/content/pdf/1297-9686-43-10.pdf>
- Ward JA, Bhangoo J, Fernandez-Fernandez F, Moore P, Swanson JD, Viola R, Velasco R, Bassil N, Weber CA, Sargent DJ (2013) Saturated linkage map construction in *Rubus idaeus* using genotyping by sequencing and genome-independent imputation. *BMC Genomics* 14(1):2. <http://www.biomedcentral.com/1471-2164/14/2>
- Young ND, DeBell F, Oldroyd GED, Geurts R, Cannon SB, Udvardi MK, Benedito VA, Mayer KFX, Gouzy J, Schoof H, Van de Peer Y, Proost S, Cook DR, Meyers BC, Spannagl M, Cheung F, De Mita S, Krishnakumar V, Gundlach H, Zhou S, Mudge J, Bharti AK, Murray JD, Naoumkina MA, Rosen B, Silverstein KAT, Tang H, Rombauts S, Zhao PX, Zhou P, Barbe V, Bardou P, Bechner M, Bellec A, Berger A, Bergs H, Bidwell S, Bisseling T, Choise N, Couloux A, Denny R, Deshpande S, Dai X, Doyle JJ, Dubez AM, Farmer AD, Fouteau S, Franken C, Gibelin C, Gish J, Goldstein S, Gonzlez AJ, Green PJ, Hallab A, Hartog M, Hua A, Humphray SJ, Jeong DH, Jing Y, Jcker A, Kenton SM, Kim DJ, Klee K, Lai H, Lang C, Lin S, Macmil SL, Magdelenat G, Matthews L, McCorrison J, Monaghan EL, Mun JH, Najjar FZ, Nicholson C, Noirot C, O'Bleness M, Paule CR, Poulain J, Prion F, Qin B, Qu C, Retzel EF, Riddle C, Sallet E, Samain S, Samson N, Sanders I, Saurat O, Scarpelli C, Schiex T, Segurens B, Severin AJ, Sherrier DJ, Shi R, Sims S, Singer SR, Sinharoy S, Sterck L, Viollet A, Wang BB, Wang K, Wang M, Wang X, Warfsmann J, Weissenbach J, White DD, White JD, Wiley GB, Wincker P, Xing Y, Yang L, Yao Z, Ying F, Zhai J, Zhou L, Zuber A, Dnari J, Dixon RA, May GD, Schwartz DC, Rogers J, Qutier F, Town CD, Roe BA (2011) The Medicago genome provides insight into the evolution of rhizobial symbioses. *Nature* 480(7378):520–524. doi:10.1038/nature10625

Terms and Conditions

Springer Nature journal content, brought to you courtesy of Springer Nature Customer Service Center GmbH (“Springer Nature”). Springer Nature supports a reasonable amount of sharing of research papers by authors, subscribers and authorised users (“Users”), for small-scale personal, non-commercial use provided that all copyright, trade and service marks and other proprietary notices are maintained. By accessing, sharing, receiving or otherwise using the Springer Nature journal content you agree to these terms of use (“Terms”). For these purposes, Springer Nature considers academic use (by researchers and students) to be non-commercial.

These Terms are supplementary and will apply in addition to any applicable website terms and conditions, a relevant site licence or a personal subscription. These Terms will prevail over any conflict or ambiguity with regards to the relevant terms, a site licence or a personal subscription (to the extent of the conflict or ambiguity only). For Creative Commons-licensed articles, the terms of the Creative Commons license used will apply.

We collect and use personal data to provide access to the Springer Nature journal content. We may also use these personal data internally within ResearchGate and Springer Nature and as agreed share it, in an anonymised way, for purposes of tracking, analysis and reporting. We will not otherwise disclose your personal data outside the ResearchGate or the Springer Nature group of companies unless we have your permission as detailed in the Privacy Policy.

While Users may use the Springer Nature journal content for small scale, personal non-commercial use, it is important to note that Users may not:

1. use such content for the purpose of providing other users with access on a regular or large scale basis or as a means to circumvent access control;
2. use such content where to do so would be considered a criminal or statutory offence in any jurisdiction, or gives rise to civil liability, or is otherwise unlawful;
3. falsely or misleadingly imply or suggest endorsement, approval, sponsorship, or association unless explicitly agreed to by Springer Nature in writing;
4. use bots or other automated methods to access the content or redirect messages
5. override any security feature or exclusionary protocol; or
6. share the content in order to create substitute for Springer Nature products or services or a systematic database of Springer Nature journal content.

In line with the restriction against commercial use, Springer Nature does not permit the creation of a product or service that creates revenue, royalties, rent or income from our content or its inclusion as part of a paid for service or for other commercial gain. Springer Nature journal content cannot be used for inter-library loans and librarians may not upload Springer Nature journal content on a large scale into their, or any other, institutional repository.

These terms of use are reviewed regularly and may be amended at any time. Springer Nature is not obligated to publish any information or content on this website and may remove it or features or functionality at our sole discretion, at any time with or without notice. Springer Nature may revoke this licence to you at any time and remove access to any copies of the Springer Nature journal content which have been saved.

To the fullest extent permitted by law, Springer Nature makes no warranties, representations or guarantees to Users, either express or implied with respect to the Springer nature journal content and all parties disclaim and waive any implied warranties or warranties imposed by law, including merchantability or fitness for any particular purpose.

Please note that these rights do not automatically extend to content, data or other material published by Springer Nature that may be licensed from third parties.

If you would like to use or distribute our Springer Nature journal content to a wider audience or on a regular basis or in any other manner not expressly permitted by these Terms, please contact Springer Nature at

onlineservice@springernature.com