

A Fast Data-Driven Method for Genotype Imputation, Phasing, and Local Ancestry Inference: MendelImpute.jl

Benjamin B. Chu¹, Eric M. Sobel^{1,3}, Rory Wasiolek¹, Janet S. Sinsheimer^{1,2,3},
Hua Zhou^{2*}, Kenneth Lange^{1,3†}

¹Department of Computational Medicine, David Geffen School of Medicine at UCLA, Los Angeles, USA

²Department of Biostatistics, Fielding School of Public Health at UCLA, Los Angeles, USA

³Department of Human Genetics, David Geffen School of Medicine at UCLA, Los Angeles, USA

keywords: Imputation; haplotyping; phasing; Julia; GWAS; admixture

1 Abstract

Current methods for genotype imputation and phasing exploit the sheer volume of data in haplotype reference panels and rely on hidden Markov models. Existing programs all have essentially the same imputation accuracy, are computationally intensive, and generally require pre-phasing the typed markers. We propose a novel data-mining method for genotype imputation and phasing that substitutes highly efficient linear algebra routines for hidden Markov model calculations. This strategy, embodied in our Julia program `MendelImpute.jl`, avoids explicit assumptions about recombination and population structure while delivering similar prediction accuracy, better memory usage, and an order of magnitude or better run-times compared to the fastest competing method. `MendelImpute` operates on both dosage data and unphased genotype data and simultaneously imputes missing genotypes and phase at both the typed and untyped SNPs. Finally, `MendelImpute` naturally extends to global and local ancestry estimation and lends itself to new strategies for data compression and hence faster data transport and sharing.

*Corresponding author. Email: huazhou@ucla.edu

†Corresponding author. Email: klange@ucla.edu

21 **2 Introduction**

22 Haplotyping (phasing) is the process of inferring unobserved haplotypes from observed genotypes. It is possi-
23 ble to deduce phase from the observed genotypes of surrounding pedigree members [25], but pedigree data are
24 no longer considered competitive with linkage disequilibrium data. Current methods for phasing and geno-
25 type imputation exploit public reference panels such as those curated by the Haplotype Reference Consortium
26 (HRC) [22] and the NHLBI TOPMed Program [27]. The sizes of these reference panels keep expanding: from
27 1000 samples in 2012 [2], to about 30,000 in 2016 [22], and to over 100,000 in 2019 [27]. Genome-wide as-
28 sociation studies (GWAS), the primary consumers of imputation, exhibit similar trends in increasing sample
29 sizes and denser SNP typing [26]. Despite these technological improvements, phasing and imputation meth-
30 ods are still largely based on hidden Markov models (HMM). Through decades of successive improvements,
31 HMM software is now more than 10,000 times faster than the original software [9], but the core HMM princi-
32 ples remain relatively unchanged. This paper explores an attractive data-driven alternative for imputation and
33 phasing that is faster and more scalable than HMM methods.

34 Hidden Markov models (HMMs) capture the linkage disequilibrium in haplotype reference panels based
35 on the probabilistic model of Li and Stephens [20]. The latest HMM software programs include Minimac
36 4 [10], Beagle 5 [6], and Impute 5 [24]. These HMMs programs all have essentially the same imputa-
37 tion accuracy [6], are computationally intensive, and generally require pre-phased genotypes. The biggest
38 computational bottleneck facing these programs is the size of the HMM state space. An initial pre-phasing
39 (imputation) step fills in missing phases and genotypes at the typed markers in a study. The easier second step
40 constructs haplotypes on the entire set of SNPs in the reference panel from the pre-phased data [13]. This
41 separation of tasks forces users to chain together different computer programs, reduces imputation accuracy
42 [9], and tends to inflate overall run-times even when the individual components are well optimized.

43 Purely data-driven techniques are potential competitors to HMMs in genotype imputation and haplotyp-
44 ing. Big data techniques substitute massive amounts of training data for detailed models in prediction. This
45 substitution can reduce computation times and, if the data are incompatible with the assumptions underly-
46 ing the HMM, improve accuracy. Haplotyping HMMs, despite their appeal and empirically satisfying error
47 rates, make simplifying assumptions about recombination hot spots and linkage patterns. We have previously
48 demonstrated the virtues of big data methods in genotype imputation with haplotyping [7] and without haplo-
49 typing [8]. SparRec [14] refines the later method by adding additional information on matrix co-clustering.
50 These two matrix completion methods efficiently impute missing entries via low-rank approximations. Un-
51 fortunately, they also rely on computationally intensive cross validation to find the optimal rank of the approx-
52 imating matrices. On the upside, matrix completion circumvents pre-phasing, exploits reference panels, and
53 readily imputes dosage data, where genotype entries span the entire interval $[0, 2]$.

54 Despite these advantages, data-driven methods have not been widely accepted as alternatives to HMM

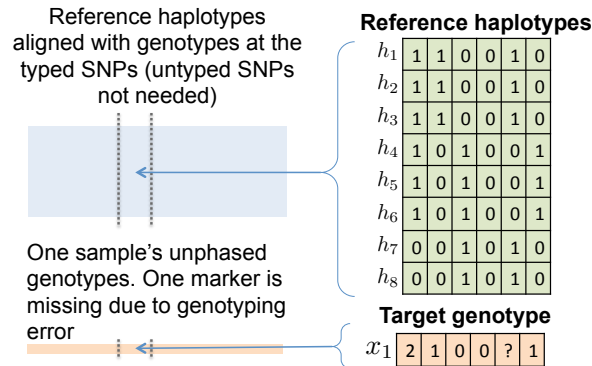
55 methods. Although it is possible in principle, our previous program [8] did not build a pipeline to handle large
56 reference panels. Here we propose a novel data-driven method to fill this gap. Our software MendelImpute
57 (a) avoids the pre-phasing step, (b) exploits known haplotype reference panels, (c) supports dosage data, (d)
58 runs extremely fast, (e) makes a relatively small demand on memory, and (f) naturally extends to local and
59 global ancestry inference. Its imputation error rate is slightly higher than the best HMM software but still
60 within a desirable range. MendelImpute is open source and forms a part of the OpenMendel platform [29],
61 which is in the modern Julia programming language [3]. We demonstrate that MendelImpute is capable of
62 dealing with HRC data even on a standard laptop. In coordination with our packages `VCFTools.jl` (handling
63 VCF files) and `SnArrays.jl` (handling PLINK files), OpenMendel powers a streamlined pipeline for end-
64 to-end data analysis. In an era where the cost of genotyping arrays continues to drop faster than Moore’s law
65 and telomere-to-telomere reference panels are within reach [23, 28], MendelImpute offers a compelling mix
66 of excellent speed, small memory footprint, and simplicity of use.

67 For each chromosome of a study subject, MendelImpute reconstructs two extended haplotypes \mathbf{E}_1 and
68 \mathbf{E}_2 that cover the entire chromosome. Both \mathbf{E}_1 and \mathbf{E}_2 are mosaics of reference haplotypes with a few break
69 points where a switch occurs from one reference haplotype to another. The break points presumably represent
70 contemporary or ancient recombination events. MendelImpute finds these reference segments and their break
71 points. From \mathbf{E}_1 and \mathbf{E}_2 it is trivial to impute missing genotypes, both typed and untyped. The extended
72 haplotypes can be painted with colors indicating the region on the globe from which each reference segment
73 was drawn. The number of SNPs assigned to each color immediately determine ethnic proportions and plays
74 into admixture mapping. The extended segments also serve as a convenient device for data compression. One
75 simply stores the break points and the index of the reference haplotype assigned to each segment. Finally, \mathbf{E}_1
76 and \mathbf{E}_2 can be nominated as maternal or paternal whenever either parent of a sample subject is also genotyped.

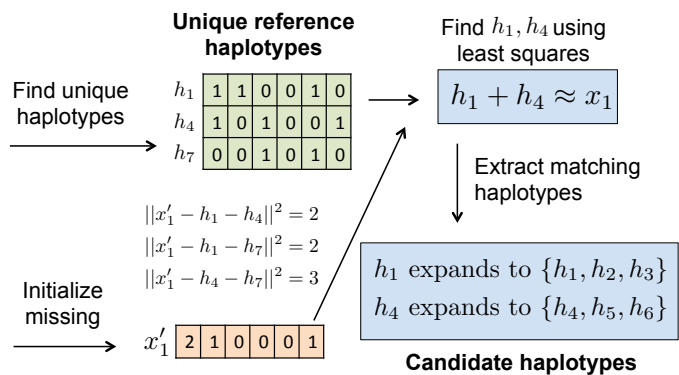
77 **3 Materials and Methods**

78 Our overall imputation strategy operates on an input matrix \mathbf{X} whose columns are sample genotypes at the
79 typed markers. The entries of \mathbf{X} represent alternative allele counts $x_{ij} \in [0, 2] \cup \{\text{missing}\}$. The reference
80 haplotypes are stored in a matrix \mathbf{H} whose columns are haplotype vectors with entries $h_{ij} \in \{0, 1\}$, representing
81 reference and alternative alleles, respectively. Given these data, the idea is to partition each sample’s genotype
82 vector into small adjacent genomic windows. In each window, many reference haplotypes collapse to the same
83 unique haplotype at the typed SNPs. We find the two unique haplotypes whose vector sum best matches the
84 sample genotype vector. Then we expand the unique haplotypes into two sets of matching full haplotypes and
85 intersect these sets across adjacent windows. Linkage disequilibrium favors long stretches of single reference
86 haplotypes punctuated by break points. Our strategy is summarized in Figure 1. A detailed commentary on
87 the interacting tactics appears in subsequent sections.

a Examine window of typed SNPs

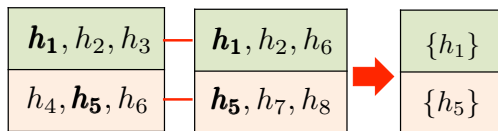


b Find optimal haplotype pairs in each window

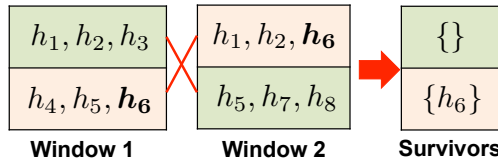


c Connect neighbors in 1 of 2 ways

Parallel connection generates 2 surviving haplotypes:

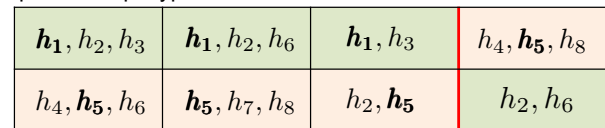


A switch at window 2 generates 1 surviving haplotype:



d Stitch window-by-window from left to right

Unphased haplotypes



Phased haplotypes

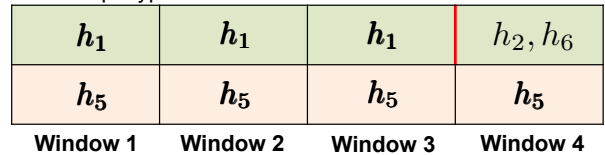


Figure 1: Overview of MendelImpute's algorithm. (a) After alignment, imputation and phasing are carried out on short, non-overlapping windows of the typed SNPs. (b) Based on a least squares criterion, we find two unique haplotypes whose vector sum approximates the genotype vector on the current window. Once this is done, all reference haplotypes corresponding to these two unique haplotypes are assembled into two sets of candidate haplotypes. (c) We intersect candidate haplotype sets window by window, carrying along the surviving set and switching orientations if the result generates more surviving haplotypes. (d) After three windows the top extended chromosome possess no surviving haplotypes, but a switch to the second orientation in the current window allows h_5 to survive on the top chromosome. Eventually we must search for a breakpoint separating h_1 from h_2 or h_6 between windows 3 and 4 (bottom panel).

88 **3.1 Missing Data in Typed and Untyped SNPs**

89 There are two kinds of missing data requiring imputation. A GWAS data set may sample as many as 10^6
90 SNPs across the genome. We call SNPs that are sampled at this stage typed SNPs. Raw data from a GWAS
91 study may contain entries missing at random due to experimental errors, but the missing rate is usually low, at
92 most a few percent, and existing programs [21] usually impute these in the pre-phasing step. When modern
93 geneticists speak of imputation, they refer to imputing phased genotypes at the unsampled SNPs present in
94 the reference panel. We call the unsampled markers untyped SNPs. The latest reference panels contain from
95 10^7 to 10^8 SNPs, so an imputation problem can have more than 90% missing data. We assume that the typed
96 SNPs sufficiently cover the entire genome. From the mosaic of typed and untyped SNPs, one can exploit local
97 linkage disequilibrium to infer for each person his/her phased genotypes at all SNPs, typed and untyped. As
98 a first step one must situate the typed SNPs among the ordered SNPs in the reference panel (Figure 1A). The
99 Julia command `indexin()` quickly finds the proper alignment.

100 **3.2 Elimination of Redundant Haplotypes by Hashing**

101 Within a small genomic window of the reference panel, multiple haplotype pairs may be identical at the
102 typed SNPs. Only the unique haplotypes play a role in matching reference haplotypes to sample genotypes.
103 `MendelImpute` identifies redundant haplotypes by hashing. For each reference haplotype limited to the win-
104 dow, hashing stores an integer representation of the haplotype via a hash function. This integer serves as an
105 index (key) to locate the reference haplotype (value). Put another way, hashing stores the inverse images of
106 the map from reference haplotypes to unique haplotypes. In our software, the `GroupSlices.jl` package [12]
107 identifies a unique key for each haplotype.

108 **3.3 Finding Optimal Haplotype Pairs via Least Squares**

109 Suppose there are d unique haplotypes $\mathbf{h}_1, \dots, \mathbf{h}_d$ (entries 0 or 1) in a genomic window (Figure 1B). Consider
110 a genotype vector \mathbf{x} with entries $x_i \in [0, 2] \cup \{\text{missing}\}$. The goal is to find the two unique haplotypes \mathbf{h}_i and
111 \mathbf{h}_j such that $\mathbf{x} \approx \mathbf{h}_i + \mathbf{h}_j$. The best haplotype pair is selected by minimizing the least squares criterion

$$\|\mathbf{x} - \mathbf{h}_i - \mathbf{h}_j\|_2^2 = \|\mathbf{x}\|_2^2 + \|\mathbf{h}_i\|_2^2 + \|\mathbf{h}_j\|_2^2 + 2\mathbf{h}_i^T \mathbf{h}_j - 2\mathbf{x}^T \mathbf{h}_i - 2\mathbf{x}^T \mathbf{h}_j$$

112 over all $\binom{d}{2} + d$ haplotype combinations. To fill in a missing value x_i , we naively initialize it with the mean
113 at each typed SNP. This action may lead to imputation errors if the typed SNPs exhibit a large proportion of
114 missingness. We discuss a strategy to remedy this bias in the supplementary sections.

115 To minimize the criterion (3.1) efficiently, suppose the genotype vectors \mathbf{x}_i constitute the columns of a
116 genotype matrix \mathbf{X} , and suppose the haplotype vectors \mathbf{h}_i constitute the columns of a haplotype matrix \mathbf{H} .

117 Given these conventions we recover all inner products $\mathbf{x}_i^T \mathbf{h}_j$ and $\mathbf{h}_i^T \mathbf{h}_j$ in equation (3.1) as entries of two
 118 matrix products; the two corresponding BLAS (Basic Linear Algebra Subroutines)[18] level-3 calls produce

$$\mathbf{X}^T \mathbf{H} = \begin{pmatrix} \mathbf{x}_1^T \mathbf{h}_1 & \cdots & \mathbf{x}_1^T \mathbf{h}_d \\ \vdots & & \vdots \\ \mathbf{x}_n^T \mathbf{h}_1 & \cdots & \mathbf{x}_n^T \mathbf{h}_d \end{pmatrix}_{n \times d} \quad \text{and} \quad \mathbf{H}^T \mathbf{H} = \begin{pmatrix} \|\mathbf{h}_1\|_2^2 & \cdots & \mathbf{h}_1^T \mathbf{h}_d \\ \vdots & & \vdots \\ \mathbf{h}_d^T \mathbf{h}_1 & \cdots & \|\mathbf{h}_d\|_2^2 \end{pmatrix}_{d \times d} .$$

119 These allow one to quickly assemble a matrix \mathbf{M} with entries $m_{ij} = \|\mathbf{h}_i\|_2^2 + \|\mathbf{h}_j\|_2^2 + 2\mathbf{h}_i^T \mathbf{h}_j$ and for each
 120 sample \mathbf{x}_k a matrix \mathbf{N} with entries $n_{ij} = -2\mathbf{x}_k^T \mathbf{h}_i - 2\mathbf{x}_k^T \mathbf{h}_j$. Therefore, to find the best haplotype pair $(\mathbf{h}_i, \mathbf{h}_j)$ for
 121 the sample \mathbf{x}_k , we search for the minimum entry $m_{ij} + n_{ij}$ of the $d \times d$ matrix $\mathbf{M} + \mathbf{N}$ across all indices $i \geq j$.
 122 Extremely unlikely ties are arbitrarily broken. Note that the term $\|\mathbf{x}_k\|_2^2$ can be safely ignored at this step.
 123 Data import and this minimum entry search are the computational bottlenecks of our software. Once such a
 124 haplotype pair is identified, all reference haplotype pairs identical to $(\mathbf{h}_i, \mathbf{h}_j)$ in the current window give the
 125 same optimal ℓ_2 error.

126 3.4 Adaptive Window Widths via Recursive Bisection

127 The width of genomic windows is an important parameter determining both imputation efficiency and accu-
 128 racy. Empirically, larger window widths give better error rates but also increase the computational burden of
 129 the matrix multiplications and minimum entry search described in Section 3.3. The magnitudes of these bur-
 130 dens depend on local haplotype diversity. Thus, we choose window widths dynamically. This goal is achieved
 131 by a bisection strategy. Initially we view all typed SNPs on a large section of a chromosome as belonging to
 132 a single window. We then divide the window into equal halves if it possesses too many unique haplotypes.
 133 Each half is further bisected and so forth recursively until every window contains fewer than a predetermined
 134 number of unique haplotypes. Empirically, choosing the maximum number d_{max} of unique haplotypes per
 135 window to be 1000 works well for both real and simulated data. When a larger number is preferred, we re-
 136 sort to a stepwise search heuristic for minimizing criterion (3.1) that scales linearly in the number of unique
 137 haplotypes d . This heuristic is described in the supplementary sections.

138 3.5 Phasing by Intersecting Haplotype Sets

139 As just described, each window w along a sample chromosome generates an optimal pair of unique haplotypes.
 140 These expand into two sets S_{w1} and S_{w2} of reference haplotypes. In the first window we arbitrarily assign S_{11}
 141 to extended haplotype 1 and S_{12} to extended haplotype 2. From here on the goal is to reconstruct two extended
 142 composite haplotypes \mathbf{E}_1 and \mathbf{E}_2 that cover the entire chromosome. Let w index the current window. The two
 143 sets $S_{w-1,1}$ and $S_{w-1,2}$ are already phased. The new sets S_{w1} and S_{w2} are not, and their phases must be resolved
 144 and their entries pruned by intersection to achieve extended haplotype parsimony. The better orientation is one

145 which generates more surviving haplotypes after intersection. Here we count surviving haplotypes across both
146 sets of an orientation. The better orientation and the corresponding survivor sets are propagated to subsequent
147 windows. If either intersection is empty at window w , then a break is declared, the empty set is replaced by
148 the entire haplotype set of window w , and a new reference segment commences. Ties and double empties
149 virtually never occur. Repeated intersection may fail to produce singleton haplotype sets, in which case we
150 randomly designate a winner to use for breakpoint search.

151 For example, suppose $S_{11} = \{\mathbf{h}_1, \mathbf{h}_2, \mathbf{h}_3\}$ and $S_{12} = \{\mathbf{h}_4, \mathbf{h}_5, \mathbf{h}_6\}$ are the (arbitrarily) phased sets in window
152 1. Since window 2 is not yet phased, the two sets $S_{21} = \{\mathbf{h}_1, \mathbf{h}_2, \mathbf{h}_6\}$ and $S_{22} = \{\mathbf{h}_5, \mathbf{h}_7, \mathbf{h}_8\}$ can be assigned
153 to extended haplotypes 1 and 2, respectively, or vice versa as depicted in Figure 1C. The first orientation
154 is preferred since it generates two surviving haplotypes \mathbf{h}_1 and \mathbf{h}_5 bridging windows 1 and 2. Thus, $\{\mathbf{h}_1\}$
155 and $\{\mathbf{h}_5\}$ are assigned at window 2 with this orientation and propagated to window 3. In window 3 the
156 contending pairs are $\{\mathbf{h}_1\} \cap \{\mathbf{h}_1, \mathbf{h}_3\}$ and $\{\mathbf{h}_5\} \cap \{\mathbf{h}_2, \mathbf{h}_5\}$ versus $\{\mathbf{h}_1\} \cap \{\mathbf{h}_2, \mathbf{h}_5\}$ and $\{\mathbf{h}_5\} \cap \{\mathbf{h}_1, \mathbf{h}_3\}$. The
157 former prevails, and $\{\mathbf{h}_1\}$ and $\{\mathbf{h}_5\}$ are assigned to window 3 and propagated to window 4. In window 4
158 the opposite orientation is preferred (Figure 1D). In this empty intersection case we set $S_{41} = \{\mathbf{h}_2, \mathbf{h}_6\}$ and
159 $S_{42} = \{\mathbf{h}_5\}$ and continue the process. Later we return and resolve the breakpoint in extended haplotype 1
160 between windows 3 and 4.

161 3.6 Resolving Breakpoints

162 The unique haplotype pairs found for adjacent windows are sometimes inconsistent and yield empty intersec-
163 tions. In such situations, we search for a good break point. Figure 1D illustrates a single-breakpoint search.
164 In this example, we slide the putative break point b across windows 3 and 4 in the top extended haplotype to
165 minimize the least squares value determined by the observed genotype, \mathbf{h}_5 spanning both windows, and the
166 breakpoint b between \mathbf{h}_1 and $\mathbf{h}_2 \cup \mathbf{h}_6$. When there is a double mismatch, we must search for a pair (b_1, b_2) of
167 breakpoints, one for each extended haplotype. The optimal pair can be determined by minimizing the least
168 squares distances generated by all possible breakpoint pairs (b_1, b_2) . Thus, double breakpoint searches scale
169 as a quadratic. Fortunately, under the adaptive window width strategy described in Section 3.4, the number
170 of typed SNPs in each window typically is on the order of 10^2 . In this range, quadratic search remains fairly
171 efficient.

172 3.7 Imputation and Phasing of Untyped SNPs

173 Once haplotyping is complete, it is trivial to impute missing SNPs. Each missing SNP is located on the
174 reference map, and its genotype is imputed as the sum of the alleles on the extended haplotypes \mathbf{E}_1 and \mathbf{E}_2 .
175 Observed genotypes are untouched unless the user prefers phased genotypes. In this case MendelImpute will
176 override observed genotypes with phased haplotypes similar to Minimac 4. Unfortunately, MendelImpute

Data Set	size (MB) for format:			
	vcf.gz	jlso	bref3	m3vcf.gz
sim 10K	48	7	18	76
sim 100K	467	54	78	33
sim 1M	4300	741	397	NA
1000G chr10	419	190	327	111
1000G chr20	183	78	148	50
HRC chr 10	3200	1101	1156	505
HRC chr 20	1500	886	529	241

Table 1: Storage size required for various compressed reference haplotype formats. Here `vcf.gz` is the standard compressed VCF format, `jlso` is used by MendelImpute, `bref3` is used by Beagle 5.1, and `m3vcf.gz` is used by Minimac 4. For all `jlso` files we chose the maximum number of unique haplotypes per window to be $d_{max} = 1000$. Note we could not generate the `m3vcf.gz` file for the sim 1M panel because it required too much memory (RAM).

177 cannot compute estimated dosages. As shown in Section 4.3 on alternative compression schemes, the extended
178 haplotypes \mathbf{E}_1 and \mathbf{E}_2 can be output rather than imputed genotypes at the user’s discretion.

179 3.8 Compressed Haplotype Panels

180 Large reference files are typically stored as compressed VCF files. Since VCF files are plain text files, they
181 are large and slow to read. Read times can be improved by computing and storing an additional tabix index
182 file [19], but file size remains a problem. Consequently, every modern imputation program has developed its
183 own specialized reference file format (for instance, the `m3vcf` and `bref3` formats of Minimac and Beagle,
184 respectively) for improving read times and storage efficiency. We propose yet another compressed format for
185 this purpose: the `jlso` format, and we compare it against other formats in Table 1.

186 The `jlso` format is constructed in three steps: (a) pre-process large reference files window-by-window,
187 (b) retain only the unique haplotypes in addition to hash maps to reference haplotypes, and (c) save the result
188 in a binary compressed format via the `JLSO.jl` package [11]. The resulting `jlso` files are 30-50x faster
189 to read and 3-5x smaller in file size (varies depending on window width) than compressed VCF files in the
190 `vcf.gz` format. Generating `jlso` files requires specifying the typed SNPs’ positions, a procedure discussed
191 in more detail in the supplementary sections. Note that in contrast to the `m3vcf` format, the `jlso` format is
192 not a standalone file format. Rather, it is a container object that facilitates reading and transferring large VCF
193 files stored as Julia variables. In principle, all files that are large in size or slow to read can be saved in this
194 alternative format for quicker access. Later we discuss an attractive alternative for storing imputed sample
195 haplotypes.

Data Set	Total SNPs	Typed SNPs	Samples	Ref Haplotypes	min MAF	Missing %
Sim 10K	62704	22879	1000	10000	0.05	0.5%
Sim 100K	80029	23594	1000	100000	0.05	0.5%
Sim 1M	97750	23092	1000	1000000	0.05	0.5%
1000G Chr10	1511445	169914	100	4808	0.25	0.1%
1000G Chr18	852602	50000*	100	4808	0.01	0.1%
1000G Chr20	669987	27091	100	4808	0.40	0.1%
HRC Chr10	1809068	116817	1000	52330	0.25	1.0%
HRC Chr20	829265	178541	1000	52330	0.01	1.0%

Table 2: Summary of real and simulated data sets. Here min MAF denotes the minimum minor allele frequency of the typed SNPs. Missing % is the percentage of typed SNPs that are randomly masked. (* We used the top 50,000 most ancestry informative SNPs.)

196 3.9 Parallel Computing and Memory Requirements

197 MendelImpute employs a shared-memory parallel computing model where each available core handles an
198 independent component of the entire problem. Work is assigned via Julia’s multi-threading functionality.
199 When computing the optimal haplotype pairs in equation (3.1), we parallelize over windows. This requires
200 allocating c copies of $\mathbf{X}^T\mathbf{H}$ and $\mathbf{H}^T\mathbf{H}$, where c is the number of CPU cores available. Note the dimensions
201 of these matrices vary across windows. To avoid accruing allocations, we pre-allocate c copies of $n \times d_{max}$
202 and $d_{max} \times d_{max}$ matrices and re-use their top-left corners in windows with $d < d_{max}$. For intersecting adjacent
203 reference haplotype sets (phasing), we parallelize over samples. This step requires no additional memory.
204 Writing to output is also trivially parallelizable by assigning each thread to write a different portion of the
205 imputed matrix to a different file, then concatenating these files into a single output file. Data import is not
206 parallelized. Beyond allocating $\mathbf{X}^T\mathbf{H}$ and $\mathbf{H}^T\mathbf{H}$, our software requires enough memory (RAM) to load the
207 target genotype matrix and the compressed haplotype reference panel.

208 3.10 Real and Simulated Data Experiments

209 For each data set, we exclude any typed SNPs with fewer than 5 copies of the minor allele and use only
210 bi-allelic SNPs. Table 2 summarizes the real and simulated data used in our comparisons.

211 3.10.1 Simulated Data

212 We simulated three 10 Mb sequence data sets, with 12,000, 102,000, 1,002,000 haplotypes, using the software
213 msprime [15]. We randomly selected 1000 samples (2000 haplotypes) from each pool to form the target
214 genotypes and used the remaining to form the reference panels. All SNPs with minor allele frequency greater
215 than 5% were designated the typed SNPs. Then 0.5% of the typed genotypes were randomly masked to

216 introduce missing values.

217 **3.10.2 1000 Genomes Data**

218 We downloaded the publicly available 1000 Genomes (1000G) phase 3 version 5a data set [1, 2]. This data set
219 contains 2504 samples with 49,143,605 phased genotypes across 26 different populations, as summarized in
220 Table 5 in the supplementary sections. We focused on chromosomes 10, 18, and 20 data in our experiments.
221 For chromosome 10 and 20, we randomly selected 100 samples to serve as imputation targets and the remain-
222 ing samples to serve as the reference panel. For chromosome 10, we chose SNPs with minor allele frequency
223 (MAF) greater than 0.25 to be typed SNPs, while for chromosome 20 we chose SNPs with MAF greater than
224 0.4 to be typed SNPs. These data sets feature in our speed and accuracy comparisons. For chromosome 18,
225 we chose the top 50,000 most ancestry informative markers (AIMs) with minor allele frequency ≥ 0.01 as the
226 typed SNPs [4]. The AIM markers were computed using `VCFTools.jl`. To avoid samples with likely substan-
227 tial degrees of continental-scale admixture, we excluded the populations ACB, ASW, CLM, GIH, ITU, MXL,
228 PEL, PUR, and STU from the reference panel. To illustrate admixture and chromosome painting, our three
229 sample individuals are taken from the excluded populations. The samples from the remaining populations are
230 assumed to exhibit less continental-scale admixture and serve as the reference panel. For admixture analysis,
231 it would be ideal to have samples from the indigenous Amerindian populations as part of the reference panel,
232 but these populations are not surveyed in the 1000 Genomes data, and so we use East Asian (EAS) and South
233 Asian (SAS) populations as the best available proxy.

234 **3.10.3 Haplotype Reference Consortium Data**

235 We also downloaded the Haplotype Reference Consortium (HRC) v1.1 data from the European Genotype-
236 Phenome archive [17] (data accession = EGAD00001002729). This data set consists of 39,741,659 SNPs in
237 27,165 individuals of predominantly European ancestry. We randomly selected 1000 samples in chromosomes
238 10 and 20 to serve as imputation targets and the remaining to serve as the reference panel. For chromosome
239 10, we selected SNPs with MAF greater than 0.25 to be typed, while for chromosome 20 we selected SNPs
240 with MAF greater than 0.01 to be typed. Finally, we randomly masked 1% of the typed genotypes to mimic
241 data missing at random.

242 **4 Results**

243 MendelImpute is publicly available at <https://github.com/OpenMendel/MendelImpute.jl>. Due to
244 Julia's flexibility, MendelImpute runs on Windows, Mac, and Linux operating systems, equipped with ei-
245 ther Intel or ARM hardware. All commands needed to reproduce the following results are available at the

246 MendelImpute site in the manuscript sub-folder.

247 **4.1 Speed, Accuracy, and Peak Memory Demand**

248 Table 3 compares the speed, accuracy, and peak memory (RAM) usage of MendelImpute, Beagle 5.1,
249 and Minimac 4. All programs were run on 10 cores of an Intel i9 9920X CPU with 64 GB of RAM. In
250 MendelImpute, the number of BLAS threads was set to 1. Note that in our simulated and real data sets the
251 correct values are known for all missing and masked genotypes. Thus, we can report accuracy as the frac-
252 tion of genotypes incorrectly imputed for all SNPs, typed or untyped. We do not compute the popular r^2
253 correlation metric for measuring imputation quality for reasons explained in the supplementary sections. All
254 reference files were previously converted to the corresponding compressed formats, bref3, m3vcf, or jls0.
255 MendelImpute was run under Julia v1.5.0. All target genotypes are unphased, and 0.1%-1% of typed geno-
256 types are missing at random. All output genotypes are phased and non-missing. For MendelImpute, we set
257 the maximum number of unique haplotypes per window to $d_{max} = 1000$. Beagle and Minimac were run under
258 their default settings. Since Minimac 4 requires pre-phased data, we used Beagle 5.1's built-in pre-phasing
259 algorithm and report its run-time and RAM usage along side those of Minimac 4.

260 On large data sets MendelImpute runs 10-30 times faster than Beagle 5.1 and 40-200 times faster than
261 Minimac 4. On the smaller 1000 Genomes data set it runs 3-5 times faster than Beagle 5.1 and 12-18 times
262 faster than Minimac 4. Increasing the reference panel size by a factor of 100 on simulated data only increases
263 MendelImpute's computation time by a factor of at most three. MendelImpute also scales better than HMM
264 methods as the number of typed SNPs increases. Thus, denser SNP arrays may benefit disproportionately
265 from using MendelImpute. The 1000 Genomes data set is exceptional in that it has fewer than 5000 reference
266 haplotypes. Therefore, traversing that HMM state space is not much slower than performing the corresponding
267 linear algebra calculations in MendelImpute. Notably, except for the HRC panels, MendelImpute spends at
268 least 50% of its total compute time importing data.

269 In terms of error rate, MendelImpute is 2-4 times worse on simulated and real data than Minimac 4
270 and Beagle 5.1. The sim10k data set is an exception in that MendelImpute's error rate is 10 times worse,
271 which we attribute to the size of the reference panel compared to the number of imputed samples. The
272 error rates of Beagle 5.1 and Minimac 4 are similar, consistent with previous findings [24]. As discussed
273 in the supplementary sections, it is possible to improve MendelImpute's error rate by more computationally
274 intensive strategies such as phasing by dynamic programming.

275 Finally, MendelImpute requires much less memory for most data sets, particularly those with a large ref-
276 erence panel or a large proportion of typed SNPs. As explained in the methods section and the supplementary
277 sections, the genotype matrix and compressed reference panel are compactly represented in memory. Since
278 most analysis is conducted in individual windows, only small sections of these matrices need to be decom-

sim 10K	Error Rate	Time (sec)	Memory (GB)
MendelImpute	3.00E-04	10	1.6
Beagle 5.1	2.81E-05	189	8.8
Minimac 4	2.38E-05	271 [177]	1.0 [6.6]
sim 100K	Error Rate	Time (sec)	Memory (GB)
MendelImpute	2.19E-05	14	1.6
Beagle 5.1	8.22E-06	279	20.1
Minimac 4	7.91E-06	3032 [253]	2.6 [14.5]
sim 1M	Error Rate	Time (sec)	Memory (GB)
MendelImpute	2.21E-05	27	4.4
Beagle 5.1	7.01E-06	769	25.6
Minimac 4	NA	NA	≥ 64
1000G chr10	Error Rate	Time (sec)	Memory (GB)
MendelImpute	1.09E-02	39	3.7
Beagle 5.1	5.51E-03	196	9.6
Minimac 4	5.24E-03	569 [159]	5.8 [10.2]
1000G chr20	Error Rate	Time (sec)	Memory (GB)
MendelImpute	3.28E-02	13	2.6
Beagle 5.1	1.68E-02	33	4.9
Minimac 4	1.65E-02	126 [33]	2.1 [5.0]
HRC chr10	Error Rate	Time (sec)	Memory (GB)
MendelImpute	6.87E-03	154	7.3
Beagle 5.1	1.90E-03	1961	32.4
Minimac 4	1.71E-03	12892 [1712]	22.5 [17.8]
HRC chr20	Error Rate	Time (sec)	Memory (GB)
MendelImpute	1.36E-03	133	6.2
Beagle 5.1	5.28E-04	2457	27.4
Minimac 4	6.34E-04	15231 [2276]	33.2 [23.5]

Table 3: Error, time, and memory comparisons on real and simulated data. Error measurement is possible in these data sets because the correct genotypes are known. Minimac 4’s benchmarks do include the pre-phasing step done by Beagle 5.1, whose memory and time are reported in brackets. For the sim 1M data, the m3vcf reference panel required for Minimac could not be computed due to excessive memory requirements. All programs were run on 10 cores of a 3.5 GHz Intel i9 CPU with 64 GB of RAM. In MendelImpute, the number of BLAS threads was set to 1.

279 pressed into single-precision arrays at any one time. Consequently, MendelImpute uses at most 7.3 GB of
280 RAM in each of these experiments. In general, MendelImpute permits standard laptops to conduct imputation
281 even with the sizeable HRC panels.

282 **4.2 Local Ancestry Inference for Admixed Populations**

283 MendelImpute lends itself to chromosome painting, the process of coloring each haplotype segment by the
284 country, region, or ethnicity of the reference individual assigned to the segment. For chromosome painting to
285 be of the most value, reference samples should be from individuals who are representative of distinct popu-
286 lations. Within a reference population there should be little admixture. Also the colors assigned to different
287 regions should be coordinated by physical, historical, and ethnic proximity. The overall proportions of the
288 colors assigned to a sample individual genome immediately translate into admixture coefficients. Here we
289 illustrate chromosome painting using chromosome 18 data from the 1000 Genomes Project. The much larger
290 Haplotype Reference Consortium data would be better suited for chromosome painting, but unfortunately its
291 repository does not list country of origin. Our examples should therefore be considered simply as a proof of
292 principle. As already mentioned, the populations present in the 1000 Genomes Project data are summarized
293 in Table 5 of the supplementary sections.

294 Figure 2A displays the painted chromosomes 18 of a native Puerto Rican (PUR, sample 1), a Peruvian
295 from Lima, Peru (PEL, sample 2), and a person of African ancestry from the Southwest USA (ASW, sample
296 3). Here a total of 17 reference populations potentially contribute genetic segments. They are colored with red,
297 brown, blue, or green to capture South Asian, East Asian, European, or African backgrounds, respectively.
298 Note that the samples from South/East Asian populations serve as a proxy for Amerindian ancestral popula-
299 tions. After coloring, the two PUR extended haplotypes are predominantly blue, the two PEL haplotypes are
300 predominantly red/brown and blue, while the two ASW haplotypes are predominantly green. Interestingly,
301 one of the PUR and one of the PEL haplotypes contain a block of African origin as well as blocks of Asian
302 and European origin, while the ASW haplotypes contain two blocks of European origin. The relatively long
303 blocks are suggestive of recent admixture. The resulting chromosome barcodes vividly display population
304 origins and suggest the locations of ancient or contemporary recombination events.

305 Figure 2B displays the three samples' admixture proportions estimated from the cumulative lengths of
306 each color in Figure 2A. From this we can immediately tell the ancestry proportion of these samples. For
307 instance, based on the chromosome 18 haplotypes, the ancestry of the Puerto Rican (sample 1; PUR) is
308 estimated to be roughly 20% S/E Asian (likely Amerindian), about 75% European, and 5% African. Of
309 course, in actual practice we would use all chromosomes of an individual, which would provide a far more
310 accurate assessment than using just chromosome 18. In addition, the ancestral assignments are only as good
311 as the choice of reference haplotypes. For instance, the haplotypes of the African American (sample 3; ASW)
312 contain a small but substantial portion (3%) of Eastern Asian ancestry. This labelling should not be taken too

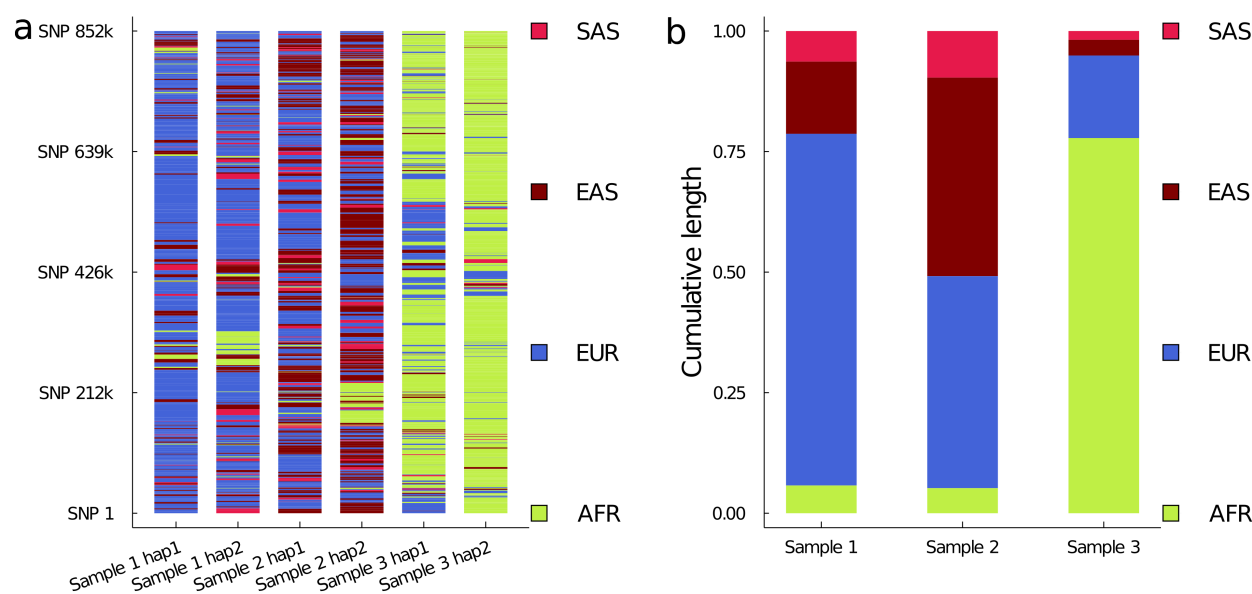


Figure 2: MendelImpute used for chromosome painting and ancestry estimation. In the reference populations, South Asians are shaded in red, East Asians in brown, Europeans in blue, and Africans in green. Sample 1 is Puerto Rican (PUR), sample 2 is Peruvian (PEL), and sample 3 is African American (ASW). (a) Painted chromosome 18 barcodes after phasing for three samples. (b) Ancestry proportions estimated from the cumulative length of the colored segments.

313 literally but rather may reflect either some distant Amerindian ancestry or regions where many haplotypes are
314 similar and so ancestry is difficult to discern.

315 4.3 Ultra-Compressed Phased Genotype Files

316 As discussed earlier, VCF files are enormous and slow to read. If genotypes are phased with respect to a
317 particular reference panel, then an alternative is to store each haplotype segment's starting position and a
318 pointer to its corresponding reference haplotype. This offers massive compression because long segments
319 are reduced to two integers. Instead of the default compressed VCF files, MendelImpute can optionally
320 output such ultra-compressed phased data. Table 4 shows that the ultra-compressed format gives 10-200 fold
321 compression compared to standard compressed VCF outputs. In principle, all phased genotypes can be stored
322 in such files. The drawback is that compressed data can only be decompressed with the help of the original
323 reference panel. Thus, this tactic relies on universal storage and curation of reference haplotype panels.
324 These panels should be stored on the cloud for easy access and constructed so that they can be consistently
325 augmented by new reference haplotypes.

Data Set	vcf.gz size (MB)	ultra-compressed size (MB)	compression ratio
Sim 10K	10.07	0.05	201
Sim 100K	10.71	0.04	267
Sim 1M	11.05	0.04	276
1000G Chr10	31.53	1.37	23
1000G Chr20	13.93	0.43	46
HRC Chr10	155.71	7.47	21
HRC Chr20	70.42	5.86	12

Table 4: Output file size comparison of compressed VCF and ultra-compressed formats.

326 4.4 Imputation Quality Scores

327 As explained in the supplementary sections, the popular r^2 correlation coefficient [5] between imputed geno-
328 types and true genotypes is an uninformative metric for measuring MendelImpute’s imputation accuracy. Al-
329 ternatively, consider the observed genotype $x_{ij} \in [0, 2] \cup \{missing\}$ at SNP i of sample j and the corresponding
330 imputed genotype g_{ij} derived from the two extended haplotypes of j . If S_i denotes the set of individuals with
331 observed genotypes at the SNP, then MendelImpute’s quality score q_i for the SNP is defined as

$$q_i = 1 - \frac{1}{|S_i|} \sum_{j \in S_i} \left(\frac{x_{ij} - g_{ij}}{2} \right)^2.$$

332 Note that $0 \leq q_i \leq 1$ and that the larger the quality score, the more confidence in the imputed values. Because
333 q_i can only be computed for the typed SNPs, an untyped SNP is assigned the average of the quality scores for
334 its two closest flanking typed SNPs. Figure 3A plots each SNP’s quality score in the 1000G Chr20 experiment
335 summarized in Table 3. For each sample, one can also compute the mean least squares error over all p SNPs to
336 obtain a per-sample quality score. This is shown in Figure 3B. By default MendelImpute outputs both quality
337 scores. Thus, investigators can perform post-imputation quality control by SNPs and by samples separately.

338 Empirically, it is rather common for a sample subject to harbor a few poorly imputed windows. Thus, we
339 observe a long right tail in the histogram for per-sample error in Figure 3. Unfortunately, the bad windows
340 generally do not exhibit any discernible regional patterns across subjects. We suspect that poorly imputed
341 windows involve breakpoints that occur near the middle of a window. We plan a detailed analysis of this issue
342 in future work.

343 5 Discussion

344 We present MendelImpute, the first scalable, data-driven method for phasing and genotype imputation.
345 MendelImpute and supporting OpenMendel software provide an end-to-end analysis pipeline in the Julia

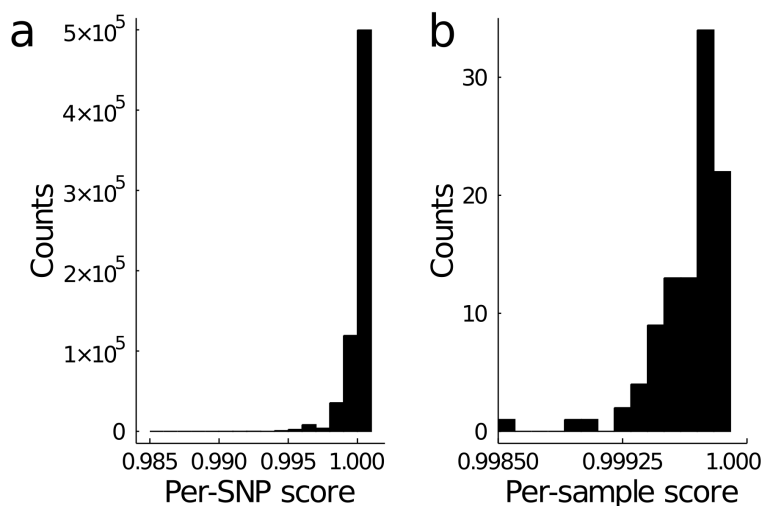


Figure 3: Histograms of per-SNP and per-sample quality scores for chromosome 20 in our 1000G analysis. By default MendelImpute computes (a) per-SNP quality scores and (b) per-sample quality scores. SNPs and samples with noticeably lower quality scores should be removed from downstream analysis.

346 programming language that is typically 10–100 times faster than methods based on hidden Markov models,
347 including Beagle 5.1 (Java) and Minimac 4 (C++). The speed difference increases dramatically as we in-
348 crease the number of typed SNPs. Thus, denser SNP chips potentially benefit more from MendelImpute’s
349 design. Furthermore, MendelImpute occupies a smaller memory footprint. This makes it possible for users
350 to run MendelImpute on standard desktop or laptop computers on very large data sets. Unfortunately we
351 cannot yet have the best of both worlds, as MendelImpute exhibits a 2-4 fold worse error rate than Beagle
352 5.1 and Minimac 4. However, as seen in Table 3, MendelImpute’s error rate is still acceptably low. One can
353 improve its error rate by implementing strategies such as phasing by dynamic programming. This strategy
354 may decrease error, but it slows down computation and is hence omitted in our comparisons. Regardless, it is
355 clear that big data methods can compete with HMM based methods on the largest data sets currently available
356 and that there is still room for improvement and innovation in genotype imputation.

357 Beyond imputation and phasing, our methods extend naturally to ancestry estimation and data compres-
358 sion. If each reference haplotype is labeled with its country or region of origin, then MendelImpute can
359 decompose a sample’s genotypes into segments of different reference haplotypes colored by these origins.
360 The cumulative lengths of these colored segments immediately yield an estimate of admixture proportions.
361 Countries can be aggregated into regions if too few reference haplotypes originate from a given country. The
362 colored segments also present a chromosome barcode that helps one visualize subject variation, recombination
363 hotspots, and global patterns of linkage disequilibrium. Data compression is achieved by storing the starting
364 positions of each segment and its underlying reference haplotype. This leads to output files that are 10-250
365 fold smaller than standard compressed VCF files. Decompression obviously requires ready access to stable

366 reference panels stored on accessible sites such as the cloud. Although such an ideal resource is currently part
367 dream and part reality, it could be achieved by a concerted international effort.

368 For potential users and developers, the primary disadvantage of MendelImpute is its reliance on the
369 importation and storage of a haplotype reference panel. Acquiring these panels requires an application process
370 which can take time to complete. Understanding, storing, and wrangling a panel add to the burden. The
371 imputation server for Minimac 4 thrives because it relieves users of these burdens [10]. Beagle 5.1 is capable
372 of fast parallel data import on raw VCF files [6], which neither Minimac 4 nor MendelImpute can currently
373 match. This makes target data import, and especially pre-processing the reference panel, painfully slow for
374 both programs. Fortunately, pre-processing only has to happen once. Beagle 5.1 also uniquely supports
375 imputation without a reference panel and directly imputes multi-allelic markers.

376 Finally, let us reiterate the goals and achievements of this paper. First, we show that data-driven meth-
377 ods are competitive with HMM methods on genotype phasing and imputation, even on the largest data sets
378 available today. Second, we challenge the notion that pre-phasing and imputation should be kept separate;
379 MendelImpute performs both simultaneously. Third, we argue that data-driven methods are ultimately more
380 flexible; for instance, MendelImpute readily handles imputation and phasing on dosage data. Fourth, we
381 demonstrate that data-driven methods yield dividends in ancestry identification and data compression. Fifth,
382 MendelImpute is completely open source, freely downloadable, and implemented in Julia, an operating sys-
383 tem agnostic, high-level programming language for scientific research. Julia is extremely fast and enables
384 clear modular coding. Our experience suggests that data-driven methods will offer a better way forward as we
385 face increasingly larger reference panels, denser SNP array chips, and greater data variability.

386 **6 Supplement**

387 **6.1 Bias Correction for Initializing Missing Data**

388 Since BLAS requires complete data, we must first initialize the missing data in each genotype vector \mathbf{x} before
389 computing \mathbf{M} and \mathbf{N} in equation (3.1). This may introduce bias in our minimization of criterion (3.1) if there is
390 a high fraction of missing genotypes in the typed SNPs, for example above 10%. One way to alleviate bias is to
391 initialize missing data with twice the alternate allele frequency and save all unique haplotype pairs minimizing
392 criterion (3.1) under this convention. Once haplotype pairs are identified, we re-minimize criterion (3.1) but
393 now skipping the missing entries of \mathbf{x} . That is equivalent to setting $x_k - h_{ik} - h_{jk} = 0$ when x_k is missing.

394 6.2 Avoidance of Global Searches for Optimal Haplotype Pairs

395 Recall that minimizing the criterion (3.1) requires searching through all lower-triangular entries of the $d \times d$
 396 matrix $\mathbf{M} + \mathbf{N}$, where d denotes the number of unique haplotypes in the window. When $d < 1000$, searching
 397 through all $\binom{d}{2} + d$ lower-triangular entries of $\mathbf{M} + \mathbf{N}$ via MendelImpute’s standard procedure is fast, but this
 398 global search quickly degrades as $d \rightarrow \infty$. Below we outline two heuristic procedures for large d .

399 6.2.1 Stepwise Search Heuristics

400 Consider minimizing the loss $f_i(\boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{x}_i - \mathbf{H}\boldsymbol{\beta}\|_2^2$, where the d columns of $\mathbf{H} \in \{0, 1\}^{p \times d}$ store unique
 401 haplotypes, p is the window width, and \mathbf{x}_i is a sample genotype vector. The original problem minimizes
 402 $f_i(\boldsymbol{\beta})$ under the constraint that exactly two $\beta_j = 1$ and the remaining $\beta_k = 0$ or the constraint that exactly one
 403 $\beta_j = 2$ and the remaining $\beta_k = 0$. As an approximate alternative, one first finds the r unique haplotypes with
 404 the largest influence on $f_i(\boldsymbol{\beta})$. This is accomplished by identifying the r most negative components of the
 405 gradient

$$\nabla f_i(\boldsymbol{\beta}) = -\mathbf{H}^T(\mathbf{x}_i - \mathbf{H}\boldsymbol{\beta}) = -\mathbf{H}^T\mathbf{x}_i + \mathbf{H}^T\mathbf{H}\boldsymbol{\beta}$$

406 at $\boldsymbol{\beta} = \mathbf{0}$. These are the r directions of steepest descent. Note that $\nabla f_i(\mathbf{0}) = -\mathbf{H}^T\mathbf{x}_i$ and that $\mathbf{H}^T\mathbf{x}_i$ is pre-
 407 computed and cached in \mathbf{N} . The residual function $g_{ij}(\mathbf{h}_k) = \frac{1}{2} \|\mathbf{x} - \mathbf{h}_j - \mathbf{h}_k\|_2^2$ is then minimized over \mathbf{h}_k to
 408 find the candidate pair $(\mathbf{h}_j, \mathbf{h}_k)$ generated by each of the vectors \mathbf{h}_j determined by the gradient $\nabla f_i(\mathbf{0})$. The
 409 ingredients to perform these minimizations are already in hand. This heuristic scales as $\mathcal{O}(rd)$, much better
 410 than $\mathcal{O}(d^2)$ in equation (3.1). MendelImpute sets the default $r = 100$. In the same spirit as the first step,
 411 one can alternatively find for each j the most negative component k of the gradient $\nabla f_i(\mathbf{e}_j)$, where \mathbf{e}_j is the
 412 standard unit vector with 1 in position j . This again determines a nearly optimal pair $(\mathbf{h}_j, \mathbf{h}_k)$. Under this
 413 tactic the Gram matrix $\mathbf{H}^T\mathbf{H}$ comes into play. Note that $\mathbf{H}^T\mathbf{H}\mathbf{e}_j$ reduces to its j th column \mathbf{v}_j . Hence, no new
 414 matrix-by-vector multiplications are necessary in calculating $\nabla f_i(\mathbf{e}_j) = -\mathbf{H}^T\mathbf{x}_i + \mathbf{v}_j = \nabla f_i(\mathbf{0}) + \mathbf{v}_j$.

415 Alternatively, one can find the best r unique haplotypes for a given sample \mathbf{x}_i *en masse* by arranging all
 416 pairwise column distances of \mathbf{X} and \mathbf{H} in the matrix

$$\mathbf{R} = \begin{bmatrix} \|\mathbf{x}_1 - \mathbf{h}_1\|_2^2 & \cdots & \|\mathbf{x}_n - \mathbf{h}_1\|_2^2 \\ \vdots & & \vdots \\ \|\mathbf{x}_1 - \mathbf{h}_d\|_2^2 & \cdots & \|\mathbf{x}_n - \mathbf{h}_d\|_2^2 \end{bmatrix}_{d \times n}$$

417 Then we partially sort each column of \mathbf{R} to identify the top r haplotypes matching each sample \mathbf{x}_i . Here \mathbf{R}
 418 is computed via the `Distance.jl` package of Julia, which internally performs BLAS level-3 calls analogous
 419 to computing $\mathbf{H}^T\mathbf{H}$ and $\mathbf{X}^T\mathbf{H}$. Instead of searching through all haplotypes to minimize $g_{ij}(\mathbf{h}_k)$ for a given
 420 sample \mathbf{x}_i , one can instead search only over the $\binom{r}{2} + r$ combinations of the top haplotypes. This allows one
 421 to entertain much larger values of r . Empirically, choosing $r = 800$ works well for most data sets.

422 **6.3 Phasing by Dynamic Programming**

423 We also investigated a dynamic programming strategy that gives the global solution for minimizing the number
424 of haplotype breaks across the extended haplotypes \mathbf{E}_1 and \mathbf{E}_2 . For each given haplotype pair $\mathbf{p}_1 = (\mathbf{h}_i, \mathbf{h}_j)$ in
425 window w , we can compute the squared Hamming distance between it and the pair $\mathbf{p}_2 = (\mathbf{h}_k, \mathbf{h}_l)$ in window
426 $w + 1$; in symbols

$$d(\mathbf{p}_1, \mathbf{p}_2) = \begin{cases} 0 & \mathbf{h}_i = \mathbf{h}_k, \mathbf{h}_j = \mathbf{h}_l & (0 \text{ breaks}) \\ 1 & \mathbf{h}_i = \mathbf{h}_k, \mathbf{h}_j \neq \mathbf{h}_l & (1 \text{ break}) \\ 1 & \mathbf{h}_i \neq \mathbf{h}_k, \mathbf{h}_j = \mathbf{h}_l & (1 \text{ break}) \\ 4 & \mathbf{h}_i \neq \mathbf{h}_k, \mathbf{h}_j \neq \mathbf{h}_l & (2 \text{ breaks}). \end{cases}$$

427 Observe that a double break is assigned an error of 4 to favor 2 single breaks across 3 windows as opposed to
428 a double break plus a perfect match.

429 Now we describe a dynamic programming strategy for finding the two paths with the minimal number
430 of unique haplotype breaks. We start with all candidate pairs \mathbf{p}_i in the leftmost window and initialize sums
431 $s_i = 0$ and traceback path vectors \mathbf{t}_i to be empty. One then recursively visits all windows in turn from left to
432 right. If w is the current window, then every candidate haplotype pair \mathbf{p}_i in window w is connected to every
433 candidate pair \mathbf{p}_j in window $w + 1$. The traceback path \mathbf{t}_j is determined by the pair \mathbf{p}_k minimizing $d(\mathbf{p}_i, \mathbf{p}_j)$.
434 The traceback path \mathbf{t}_j is constructed by appending \mathbf{p}_k to \mathbf{t}_k and setting $s_j = s_k + d(\mathbf{p}_k, \mathbf{p}_j)$. This process is
435 continued until the rightmost window v is reached. At this point the pair \mathbf{p}_j with lowest running sum s_j is
436 declared the winner. The traceback path \mathbf{t}_j allows one to construct the extended haplotypes \mathbf{E}_1 and \mathbf{E}_2 in their
437 entirety. Unfortunately, too many haplotype pairs per window can overwhelm dynamic programming with
438 large reference panels. One partial recourse is to discard partial paths and associated partial termini \mathbf{p}_i that are
439 unpromising in the sense that their running sums s_i are excessively large. This is effective, but the approximate
440 algorithm is burdened by extra bookkeeping. The bookkeeping of the exact algorithm is already demanding.

441 **6.4 Generating JLSO Compressed Reference Haplotype Panels**

442 Since haplotypes are long binary vectors, the entire haplotype reference panel can be compactly represented
443 using a single bit per entry. As a first step in compression, the reference panel is divided into non-overlapping
444 windows of various widths after proper alignment of all typed and reference SNPs in the window. For each
445 window we save two compressed mini-panels. The first houses the unique haplotypes determined by just
446 the typed SNPs. The second houses the unique haplotypes determined by all SNPs in the window, typed or
447 untyped. The former is much smaller than the later, but each entry of both can be compactly represented by a
448 single bit per entry in memory. Thus, for each window we save two compressed windows in addition to meta
449 information and pointers that coordinate reference haplotypes with the two mini-panels per each window.

450 MendelImpute also requires the list of typed SNPs. This whole ensemble is stored in the jlso file. Because
451 there are only a limited number of SNP array chips on the market, one can in principle store just a few jlso
452 files on a universal source such as the cloud.

453 6.5 Imputation quality scores using r-squared

454 The popular r^2 correlation coefficient between imputed genotypes and true genotypes is a popular metric
455 for measuring imputation quality. Although the true genotypes are unknown, r^2 can still be estimated from
456 posterior probabilities of the HMM model [5]. For instance, according to Minimac 3's online documentation,
457 for each SNP we can calculate

$$r^2 = \frac{\frac{1}{2n} \sum_{i=1}^{2n} (D_i - \hat{p})^2}{\hat{p}(1 - \hat{p})}$$

458 where \hat{p} is the alternative allele frequency (of the imputed data), D_i is the imputed alternate allele probability
459 at the i th haplotype, and n is the number of GWAS samples. MendelImpute does not compute posterior
460 probabilities; for each locus, each haplotype is imputed with 0 or 1. That is, $D_i \in \{0, 1\}$. Thus, using the fact
461 that $\hat{p} = \frac{1}{2n} \sum_{i=1}^{2n} D_i$, we have

$$\begin{aligned} r^2 &= \frac{\frac{1}{2n} \sum_{i=1}^{2n} (D_i^2 - 2D_i\hat{p} + \hat{p}^2)}{\hat{p}(1 - \hat{p})} \\ &= \frac{\frac{1}{2n} [\sum_{i=1}^{2n} D_i^2 - 2\hat{p} \sum_{i=1}^{2n} D_i + 2n\hat{p}^2]}{\hat{p}(1 - \hat{p})} \\ &= \frac{\frac{1}{2n} [\sum_{i=1}^{2n} D_i - 4n\hat{p}^2 + 2n\hat{p}^2]}{\hat{p}(1 - \hat{p})} \\ &= \frac{\frac{1}{2n} [2n\hat{p} - 2n\hat{p}^2]}{\hat{p}(1 - \hat{p})} \\ &= 1 \end{aligned}$$

462 Thus every SNP has $r^2 = 1$ under MendelImpute's model, which renders it meaningless.

463 6.6 Summary of 1000 Genomes Reference Panel

464 A total of 26 different populations contribute to the 1000 Genomes Project data set. These populations are
465 further organized into five super population. While this information is freely available online, we summarize
466 it in Table 5 for completeness.

Population Code	Description	Super Population
CHB	Han Chinese in Beijing, China	East Asian (EAS)
JPT	Japanese in Tokyo, Japan	East Asian (EAS)
CHS	Southern Han Chinese	East Asian (EAS)
CDX	Chinese Dai in Xishuangbanna, China	East Asian (EAS)
KHV	Kinh in Ho Chi Minh City, Vietnam	East Asian (EAS)
CEU	Utah Residents with NW European Ancestry	European (EUR)
TSI	Toscans in Italia	European (EUR)
FIN	Finnish in Finland	European (EUR)
GBR	British in England and Scotland	European (EUR)
IBS	Iberian Population in Spain	European (EUR)
YRI	Yoruba in Ibadan, Nigeria	Africans (AFR)
LWK	Luhya in Webuye, Kenya	Africans (AFR)
GWD	Gambian in Western Divisions in the Gambia	Africans (AFR)
MSL	Mende in Sierra Leone	Africans (AFR)
ESN	Esan in Nigeria	Africans (AFR)
ASW	Americans of African Ancestry in SW USA	Africans (AFR)
ACB	African Caribbeans in Barbados	Africans (AFR)
MXL	Mexican Ancestry from Los Angeles USA	Ad Mixed American (AMR)
PUR	Puerto Ricans from Puerto Rico	Ad Mixed American (AMR)
CLM	Colombians from Medellin, Colombia	Ad Mixed American (AMR)
PEL	Peruvians from Lima, Peru	Ad Mixed American (AMR)
GIH	Gujarati Indian from Houston, Texas	South Asian (SAS)
PJL	Punjabi from Lahore, Pakistan	South Asian (SAS)
BEB	Bengali from Bangladesh	South Asian (SAS)
STU	Sri Lankan Tamil from the UK	South Asian (SAS)
ITU	Indian Telugu from the UK	South Asian (SAS)

Table 5: The 26 population codes present in the 1000 genomes project.

467 **7 Availability of source code**

468 **Project name:** MendelImpute.jl

469 **Project home page:** <https://github.com/OpenMendel/MendelImpute.jl>

470 **Supported operating systems:** Mac OS, Linux, Windows

471 **Programming language:** Julia 1.5

472 **License:** MIT

473 **8 Author contributions**

474 KL, JS, ES, and HZ conceived this project. BC, KL, RW, JS, ES, and HZ devised the methods. BC, RW,
475 and HZ developed the software. BC and KL accessed the data. BC wrote the original draft of the paper. BC,
476 KL, RW, JS, ES, and HZ reviewed and edited the draft. KL previewed some of the methods incorporated in
477 MendelImpute in his talk [16].

478 **9 Competing interests**

479 The authors declare no competing interests.

480 **10 Acknowledgements**

481 BC and RW were supported by NIH grant T32-HG002536. BC, ES, JS, HZ, and KL were supported by
482 NIH grant R01-HG006139. ES, JS, HZ, and KL were supported by NIH grant R01-GM053275. JS was also
483 supported by NIH grant R01-HG009120.

484 We would also like to thank Calvin Chi for his helpful discussions on ancestry estimation and Juhyun Kim
485 for her helpful discussions on imputation quality scores.

486 **References**

- 487 [1] 1000 Genomes Project Consortium, A. Auton, G. R. Abecasis, D. M. Altshuler, R. M. Durbin, et al. A
488 global reference for human genetic variation. *Nature*, 526(7571):68–74, 2015.
- 489 [2] 1000 Genomes Project Consortium, G. A. McVean, D. M. Altshuler, R. M. Durbin, et al. An integrated
490 map of genetic variation from 1,092 human genomes. *Nature*, 491(7422):56–65, 2012.
- 491 [3] J. Bezanson, A. Edelman, S. Karpinski, and V. B. Shah. Julia: A fresh approach to numerical computing.
492 *SIAM Review*, 59:65–98, 2017.
- 493 [4] R. Brown and B. Pasaniuc. Enhanced methods for local ancestry assignment in sequenced admixed
494 individuals. *PLoS Comput Biol*, 10(4):e1003555, 2014.
- 495 [5] B. L. Browning and S. R. Browning. A unified approach to genotype imputation and haplotype-phase
496 inference for large data sets of trios and unrelated individuals. *The American Journal of Human Genetics*,
497 84(2):210–223, 2009.
- 498 [6] B. L. Browning, Y. Zhou, and S. R. Browning. A one-penny imputed genome from next-generation
499 reference panels. *The American Journal of Human Genetics*, 103(3):338–348, 2018.
- 500 [7] G. K. Chen, K. Wang, A. H. Stram, E. M. Sobel, and K. Lange. Mendel-gpu: haplotyping and genotype
501 imputation on graphics processing units. *Bioinformatics*, 28(22):2979–2980, 2012.
- 502 [8] E. C. Chi, H. Zhou, G. K. Chen, D. O. Del Vecchio, and K. Lange. Genotype imputation via matrix
503 completion. *Genome Research*, 23(3):509–518, 2013.
- 504 [9] S. Das, G. R. Abecasis, and B. L. Browning. Genotype imputation from large reference panels. *Annual*
505 *Review of Genomics and Human Genetics*, 19:73–96, 2018.
- 506 [10] S. Das, L. Forer, S. Schönherr, C. Sidore, A. E. Locke, A. Kwong, S. I. Vrieze, E. Y. Chew, S. Levy,
507 M. McGue, et al. Next-generation genotype imputation service and methods. *Nature Genetics*,
508 48(10):1284, 2016.
- 509 [11] R. Finnegan and L. White. invenia/JLSO.jl: Storage container for serialized Julia objects. <https://doi.org/10.5281/zenodo.3992374>, 2020.
510
- 511 [12] A. GreenWell and M. Abbott. GroupSlices.jl: A package for the groupslices and associated functions.
512 <https://github.com/mcabbott/GroupSlices.jl>, 2019.
- 513 [13] B. Howie, C. Fuchsberger, M. Stephens, J. Marchini, and G. R. Abecasis. Fast and accurate genotype
514 imputation in genome-wide association studies through pre-phasing. *Nature Genetics*, 44(8):955–959,
515 2012.
- 516 [14] B. Jiang, S. Ma, J. Causey, L. Qiao, M. P. Hardin, I. Bitts, D. Johnson, S. Zhang, and X. Huang. SparRec:
517 An effective matrix completion framework of missing data imputation for GWAS. *Scientific Reports*,
518 6:35534, 2016.
- 519 [15] J. Kelleher, A. M. Etheridge, and G. McVean. Efficient coalescent simulation and genealogical analysis
520 for large sample sizes. *PLoS Comput Biol*, 12(5):1–22, 05 2016.

- 521 [16] K. Lange. Lecture on Ultrafast Haplotyping. In *New Statistical Methods for Family-Based*
522 *Sequencing Studies*. Banff International Research Station, [http://www.birs.ca/events/2018/](http://www.birs.ca/events/2018/5-day-workshops/18w5154/videos/watch/201808091354-Lange.html)
523 [5-day-workshops/18w5154/videos/watch/201808091354-Lange.html](http://www.birs.ca/events/2018/5-day-workshops/18w5154/videos/watch/201808091354-Lange.html), 2018.
- 524 [17] I. Lappalainen, J. Almeida-King, V. Kumanduri, A. Senf, J. D. Spalding, G. Saunders, J. Kandasamy,
525 M. Caccamo, R. Leinonen, B. Vaughan, et al. The European genome-phenome archive of human data
526 consented for biomedical research. *Nature Genetics*, 47(7):692–695, 2015.
- 527 [18] C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh. Basic linear algebra subprograms for
528 Fortran usage. *ACM Transactions on Mathematical Software*, 5(3):308–323, 1979.
- 529 [19] H. Li. Tabix: fast retrieval of sequence features from generic tab-delimited files. *Bioinformatics*,
530 27(5):718–719, 2011.
- 531 [20] N. Li and M. Stephens. Modeling linkage disequilibrium and identifying recombination hotspots using
532 single-nucleotide polymorphism data. *Genetics*, 165(4):2213–2233, 2003.
- 533 [21] P.-R. Loh, P. Danecek, P. F. Palamara, C. Fuchsberger, Y. A. Reshef, H. K. Finucane, S. Schoenherr,
534 L. Forer, S. McCarthy, G. R. Abecasis, et al. Reference-based phasing using the haplotype reference
535 consortium panel. *Nature Genetics*, 48(11):1443, 2016.
- 536 [22] S. McCarthy, S. Das, W. Kretzschmar, O. Delaneau, A. R. Wood, A. Teumer, H. M. Kang, C. Fuchs-
537 berger, P. Danecek, K. Sharp, et al. A reference panel of 64,976 haplotypes for genotype imputation.
538 *Nature Genetics*, 48(10):1279, 2016.
- 539 [23] K. H. Miga, S. Koren, A. Rhie, M. R. Vollger, A. Gershman, A. Bzikadze, S. Brooks, E. Howe, D. Porub-
540 sky, G. A. Logsdon, et al. Telomere-to-telomere assembly of a complete human X chromosome. *bioRxiv*,
541 735928, 2019.
- 542 [24] S. Rubinacci, O. Delaneau, and J. Marchini. Genotype imputation using the positional burrows wheeler
543 transform. *bioRxiv*, 797944, 2019.
- 544 [25] E. Sobel, K. Lange, J. R. O’Connell, and D. E. Weeks. Haplotyping algorithms. In *Genetic Mapping*
545 *and DNA Sequencing*, pages 89–110. Springer, 1996.
- 546 [26] C. Sudlow, J. Gallacher, N. Allen, V. Beral, P. Burton, J. Danesh, P. Downey, P. Elliott, J. Green, M. Lan-
547 dray, B. Liu, P. Matthews, G. Ong, J. Pell, A. Silman, A. Young, T. Sprosen, T. Peakman, and R. Collins.
548 UK BioBank: an open access resource for identifying the causes of a wide range of complex diseases of
549 middle and old age. *PLoS Medicine*, 12:e1001779, 2015.
- 550 [27] D. Taliun, D. N. Harris, M. D. Kessler, J. Carlson, Z. A. Szpiech, R. Torres, S. A. G. Taliun, A. Corvelo,
551 S. M. Gogarten, H. M. Kang, et al. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed
552 Program. *bioRxiv*, 563866, 2019.
- 553 [28] S.-K. Yoo, C.-U. Kim, H. L. Kim, S. Kim, J.-Y. Shin, N. Kim, J. S. W. Yang, K.-W. Lo, B. Cho,
554 F. Matsuda, et al. NARD: whole-genome reference panel of 1779 Northeast Asians improves imputation
555 accuracy of rare and low-frequency variants. *Genome Medicine*, 11(1):1–10, 2019.
- 556 [29] H. Zhou, J. S. Sinsheimer, D. M. Bates, B. B. Chu, C. A. German, S. S. Ji, K. L. Keys, J. Kim, S. Ko,
557 G. D. Mosher, J. C. Papp, E. M. Sobel, J. Zhai, J. J. Zhou, and K. Lange. OpenMendel: a cooperative
558 programming project for statistical genetics. *Human Genetics*, 139:61–71, 2020.