# Choosing subsamples for sequencing studies by minimizing the average distance to the closest leaf

Jonathan T. L. Kang[*], Peng Zhang[§], Sebastian Zöllner[†], and Noah A. Rosenberg[*]

[*]Department of Biology, Stanford University, Stanford, CA 94305
[§]Center for Inherited Disease Research, Johns Hopkins University, Baltimore, MD 21224
[†]Department of Biostatistics, and Department of Psychiatry, University of Michigan, Ann Arbor, MI 48109

April 7, 2015

11 **Running title:** Internal reference panels for imputation

12 **Key words:** algorithms, imputation, polymorphic sites, sequencing, study design

13 **Corresponding author:**

14 Jonathan T. L. Kang

15 Department of Biology, Stanford University

16 371 Serra Mall, Stanford, CA 94305

17 (650) 724-5122

18 jtlkang@stanford.edu

## Abstract

19

20    Imputation of genotypes in a study sample can make use of sequenced or densely genotyped

21    external reference panels consisting of individuals that are not from the study sample. It can

22    also employ internal reference panels, incorporating a subset of individuals from the study

23    sample itself. Internal panels offer an advantage over external panels, as they can reduce

24    imputation errors arising from genetic dissimilarity between a population of interest and a

25    second, distinct population from which the external reference panel has been constructed.

26    As the cost of next-generation sequencing decreases, internal reference panel selection is

27    becoming increasingly feasible. However, it is not clear how best to select individuals to

28    include in such panels. We introduce a new method for selecting an internal reference panel—

29    minimizing the average distance to the closest leaf (ADCL)—and compare its performance

30    relative to an earlier algorithm: maximizing phylogenetic diversity (PD). Employing both

31    simulated data and sequences from the 1000 Genomes Project, we show that ADCL provides

32    a significant improvement in imputation accuracy, especially for imputation of sites with low-

33    frequency alleles. This improvement in imputation accuracy is robust to changes in reference

34    panel size, marker density, and length of the imputation target region.

## Introduction

35

36    Owing to the existence of genetic variation within species, geneticists routinely make choices

37    about which individuals, inbred strains, or representatives of populations or breeds merit

38    prioritization for genotyping or DNA sequencing. Often, such choices, though typically

39    made by informal criteria, reflect an explicit or implicit goal of maximizing the potential for

40    extrapolating the information in the genotyped or sequenced individuals to all members of

41    a breed, population, or species of interest.

Genotype imputation algorithms infer unobserved genotypes by matching a set of markers to the haplotype patterns observed in a reference sample (LI *et al.* 2009; MARCHINI AND HOWIE 2010), adding a new dimension to these choices. Reference panels that are used to facilitate genotype imputation in other individuals beyond the members of the panels themselves can often be optimally selected to formally maximize the imputed genotypic information obtained about those other individuals of interest (KANG AND MARJORAM 2012; ZHANG *et al.* 2013; PEIL *et al.* 2015). The evaluation of alternative ways to select imputation reference panels thus provides an approach for making sample choices for major genotyping or sequencing studies more systematically generalizable.

When conducting genotype imputation studies in a population sample, reference panels have generally been selected from databases external to the sample, such as the 1000 Genomes Project (1000 GENOMES PROJECT CONSORTIUM 2010) and the International HapMap Consortium (INTERNATIONAL HAPMAP CONSORTIUM 2005) databases. As a result of the rapidly decreasing cost of sequencing, however, it has become increasingly possible to carry out *internal* reference panel selection, in which additional sequencing is performed on a subset of the study sample, and the sequenced subset is then used to impute the remaining haplotypes. The use of reference sequences that originate from the study sample itself can reduce the potential mismatch of ancestral backgrounds between sample and reference populations, decreasing imputation errors. It also allows for genetic variants unique to the sample population to be successfully imputed (FRIDLEY *et al.* 2010; ZHANG *et al.* 2013).

Previous studies have observed that a mismatch in population origins between reference panels and study samples can reduce imputation accuracy compared to when they originate from the same or similar populations (HUANG *et al.* 2009, 2011; LI *et al.* 2010; PAŞANIUC *et al.* 2010; SHRINER *et al.* 2010; SURAKKA *et al.* 2010). JEWETT *et al.* (2012) demonstrated, using a coalescent model, that with other variables held constant, smaller internal reference panels are often likely to outperform larger external reference panels, despite the difference

4

<sub>68</sub> in panel size. Empirical studies have also shown that using an internal reference panel drawn

<sub>69</sub> from a subset of the sample under study, in addition to an external reference panel, gives

<sub>70</sub> rise to an increase in imputation accuracy over just using the external reference panel alone

<sub>71</sub> (FRIDLEY *et al.* 2010; SAMPSON *et al.* 2012; KREINER-MØLLER *et al.* 2015).

<sub>72</sub> The value of internal reference panels for imputation studies raises the question of how an in-

<sub>73</sub> ternal panel should be selected. Two recent studies have proposed maximizing "phylogenetic

<sub>74</sub> diversity" (PD) as a criterion for internal reference panel selection (KANG AND MARJORAM

<sub>75</sub> 2012; ZHANG *et al.* 2013). In this approach, the phylogenetic diversity of a set of haplotypes

<sub>76</sub> is defined as the total branch length of a tree spanned by the haplotypes (FAITH 1992; HART-

<sub>77</sub> MANN AND STEEL 2007). Given a panel size, the goal is to select the subset of haplotypes

<sub>78</sub> whose subtree yields the longest total branch length. Conceptually, the idea of seeking a

<sub>79</sub> maximally diverse subset of haplotypes in the reference panel aims to sample haplotypes that

<sub>80</sub> best cover the full range of haplotypes observed in the sample. The maximum-PD panel, by

<sub>81</sub> choosing haplotypes from different regions of the tree of haplotypes (FIGURE 1B), is more

<sub>82</sub> likely than a random panel to supply the necessary diversity to impute sites localized in a

<sub>83</sub> subgroup within the entire sample population. ZHANG *et al.* (2013) showed, using simulated

<sub>84</sub> sequence data and data from the 1000 Genomes Project, that by using the maximum-PD

<sub>85</sub> panel, higher imputation accuracy is obtained, and more sites are imputed as polymorphic in

<sub>86</sub> the sample population, than if the reference panel consists of randomly-selected haplotypes.

<sub>87</sub> Despite the utility of maximizing PD as a method for the selection of an internal reference

<sub>88</sub> panel, other approaches focusing on different principles might be preferable. Because the

<sub>89</sub> algorithm explicitly chooses haplotypes that are genetically distant from one another, long,

<sub>90</sub> pendant branches of the tree, if present, are likely to be chosen (BORDEWICH *et al.* 2008).

<sub>91</sub> The haplotypes associated with such branches might not be representative of the sample at

<sub>92</sub> large. These haplotypes might contain a large amount of sequencing error or missing data,

<sub>93</sub> and their inclusion in the reference panel might not contribute substantially to an increase

5

94 in imputation accuracy. Even if they have high-quality data, such haplotypes are relatively

95 unique in the sample, and therefore might assist as imputation templates only for a small

96 number of sampled lineages.

97 PD can be viewed as emphasizing "diversity" of the internal reference panel rather than

98 "representativeness." To determine if an alternative focused on identifying the most repre-

99 sentative subsample for use as the internal reference panel is preferable, we explore a new

100 method: minimizing the average distance to the closest leaf (ADCL), which identifies refer-

101 ence haplotypes based on their genetic proximity to the rest of the sample haplotypes. We

102 compare the imputation accuracy of the maximum-PD, minimum-ADCL, and random refer-

103 ence panels on both simulated data and data from the 1000 Genomes Project, and find that

104 the minimum-ADCL panel consistently provides higher imputation accuracy, irrespective of

105 changes to parameters such as reference panel size, marker density, and sequence length.

## METHODS

### Maximizing phylogenetic diversity (PD)

108 Given a tree of $n$ haplotypes, to select a reference panel of haplotypes whose subtree spans

109 the longest branch length, ZHANG *et al.* (2013) considered a greedy algorithm that takes as

110 inputs the tree and a parameter $k \leq n$, the desired number of haplotypes for the panel. Let

111 $X$ be the $k$-element subset of the sample haplotypes chosen for the reference panel, and let

112 $T_X$ be the subtree spanned by the haplotypes in $X$. The algorithm first selects the haplotype

113 pair that is phylogenetically most distant (i.e. largest pairwise branch length), and adds both

114 haplotypes to $X$. $T_X$ now consists of a single pair of branches. Sequentially, the haplotype

115 that is the most distant from $T_X$ is placed into $X$, updating $T_X$ with each inclusion. This

116 process continues until the required $k$ haplotypes have been selected (FIGURE 1B).

117   PARDI AND GOLDMAN (2005) and STEEL (2005) proved that among all possible subsets of

118   size $k \leq n$ haplotypes from the study sample, the greedy algorithm achieves the globally

119   maximal PD. Thus, the selection of the "most diverse" reference panel is computationally

120   efficient, as there is no need to exhaustively examine all possible panels of size $k$ in order

121   to arrive at the correct solution. In addition, because the selection algorithm is greedy, the

122   haplotypes in the reference panel can be ranked by their order of inclusion, in which every

123   haplotype added contributes a non-increasing amount of PD. The maximum-PD panels of

124   size 2 to $k$ form a series of nested sets, and all previously selected haplotypes in a panel of

125   size smaller than $k$ will also be included in a panel of size $k$.

### Minimizing the average distance to the closest leaf (ADCL)

127   **Overview of ADCL:** Instead of focusing on diversity in the selected set and targeting

128   the potential for accurate imputation of unusual haplotypes, the minimum-ADCL algorithm

129   focuses on representativeness, aiming to maximize imputation accuracy of typical haplotypes

130   likely to appear in a sample. The problem can be viewed as choosing the haplotypes that are,

131   on average, genealogically closest to the remaining haplotypes not included in the reference

132   panel. As in the case of PD, the algorithm takes as inputs a tree of the $n$ haplotypes in the

133   study sample, and a parameter $k \leq n$, indicating the desired reference panel size.

134   Let $H$ be the set of $n$ haplotypes, and let $X$ be the selected $k$-element subset of $H$. The

135   objective is then to find $X$ such that the branch-length distance from a randomly-chosen

136   haplotype in $H$ to its closest neighboring haplotype in $X$ is minimized over all possible $k$-

137   element subsets of $H$ (MATSEN *et al.* 2013). Note that because the haplotypes in $X$ are also

138   in $H$, each of these haplotypes is its own closest neighbor, and we can equivalently consider

139   either $H$ or $H \setminus X$. In essence, the goal is to return a set of reference panel haplotypes that

140   occupy the most central positions within clusters of the tree (FIGURE 1C).

141 In a detailed study of ADCL, MATSEN *et al.* (2013) demonstrated that unlike when choosing

142 the subset that maximizes PD, the greedy algorithm need not give rise to the globally-optimal

143 ADCL solution. It is therefore necessary to produce alternative algorithms that seek to

144 minimize ADCL. Note that because the greedy algorithm is not applicable, the haplotypes

145 selected cannot be ranked by their order of inclusion, as a haplotype included in a subset of

146 size smaller than $k$ is not necessarily also included in a subset of size $k$ (FIGURE 1C).

**Adapted partitioning-around-medoids (PAM) algorithm for minimizing ADCL:**

148 MATSEN *et al.* (2013) described two algorithms which, for a given set of haplotypes, seek to

149 produce the subset of size $k$ that minimizes ADCL. The first approach leverages similarities

150 between the problem of minimizing ADCL and the technique known as $k$-medoids clustering

151 (KAUFMAN AND ROUSSEEUW 1987). In the $k$-medoids problem, a set of data points is

152 partitioned into $k$ clusters, where $k$ is predetermined. Within each cluster, a single point is

153 designated as the center. The $k$-medoids clustering method is similar to $k$-means clustering.

154 In the $k$-medoids approach, however, each cluster center is chosen from the original set of data

155 points, whereas $k$-means has no such restriction. The objective function to be minimized in

156 the $k$-medoids problem is the distance from a random data point to the center of the cluster

157 to which it is assigned. A cluster center can be viewed as the data point most representative

158 of the remainder of the data points within the cluster.

159 It is then clear how the problem of minimizing ADCL is analogous to the $k$-medoids problem.

160 A data point is a haplotype, and distances between data points are branch-length (patristic)

161 distances between haplotypes. The $k$ cluster centers are akin to the $k$ haplotypes that are

162 selected.

163 As with minimizing ADCL, there is no greedy algorithm that solves the $k$-medoids prob-

164 lem, and obtaining the globally optimal solution has been demonstrated to be NP-hard

165 (SHENG AND LIU 2004). A widely-used $k$-medoids heuristic algorithm is the partitioning-

8

166 around-medoids (PAM) algorithm (THEODORIDIS AND KOUTROUMBAS 2008), which works

167 by randomly selecting $k$ medoids from the original set of $n$ data points, and then minimizing

168 the objective function via hill-climbing. One iteration of the algorithm consists of looping

169 over all $k(n-k)$ possible pairs containing a medoid and non-medoid, exchanging the medoid

170 statuses of the points in the pair, and recording the new value of the objective function from

171 the updated arrangement. Among all $k(n-k)$ proposed exchanges, the single exchange that

172 leads to the lowest-cost configuration is chosen. The algorithm then enters a new iteration,

173 and the process repeats until no further changes to the set of medoids take place.

174 The first approach MATSEN *et al.* (2013) considered for minimizing ADCL is an adaptation

175 of the PAM algorithm. First, the set $X$ of haplotypes included in the reference panel is

176 initialized by randomly selecting, without replacement, $k$ haplotypes from the initial set $H$

177 of $n$ haplotypes. Next, the following loop over the haplotypes $x_1, \ldots, x_k \in X$ is executed

178 until no exchanges occur for one complete iteration over every $x_i \in X$:

179 (1) For a haplotype $x_i \in X$, remove it from $X$ and attempt to replace it with every other

180      $y \in H \setminus X$ in its place.

181 (2) Keep the best such exchange if it decreases ADCL.

182 (3) Continue with $x_{i+1} \in X$. In the case of $x_k$, continue with $x_1$.

183 This method for minimizing ADCL differs from the original formulation of the PAM algo-

184 rithm in that it evaluates potential exchanges one medoid at a time, instead of examining

185 all $k(n-k)$ medoid/non-medoid pairs before finding the exchange that most decreases the

186 objective function (MATSEN *et al.* 2013). Because each step in the iteration causes the value

187 of ADCL to either stay constant or decrease, the solution is guaranteed to converge on a

188 local minimum. However, the algorithm remains a heuristic approach, and the minimum-

189 ADCL solution it achieves could depend on the specific haplotypes selected during random

190 initialization. Hence, the global minimum might not always be found.

9

191 Alongside the adapted PAM algorithm, MATSEN *et al.* (2013) also developed a second ap-

192 proach: an exact but more computationally-intensive algorithm that is guaranteed to find

193 the global-minimum ADCL solution. Both algorithms were implemented in the `rppr` bi-

194 nary in the `pplacer` suite of programs. Comparing between the two, MATSEN *et al.* (2013)

195 demonstrated that for their simulated test sets, the adapted PAM algorithm only rarely gets

196 trapped in local minima. For computational efficiency, we therefore chose to use the adapted

197 PAM algorithm rather than the slower exact algorithm, first testing that in our setting, mul-

198 tiple runs of the adapted PAM algorithm with different initial seeds select a large percentage

199 of the same haplotypes (see RESULTS).

## Simulated sequence data

201 To evaluate how the maximum-PD and minimum-ADCL panels perform relative to one an-

202 other, we analyzed simulated data sets produced by the coalescent-based sequence sampling

203 program `ms` (HUDSON 2002), closely following the parameters used by ZHANG *et al.* (2013)

204 to ensure that the results are comparable.

205 First, we independently generated 50 data sets, each consisting of 2000 1Mb haplotypes,

206 assuming a constant effective population size of $N_e = 10,000$, a mutation rate of $\mu = 10^{-8}$

207 per site per generation, and a recombination rate of $\rho = 10^{-8}$ per site per generation. The

208 parameter values provided to `ms` were as follows: $\mathtt{nsam} = 2000$, $\mathtt{nreps} = 50$, $\mathtt{-t} = 400$, $\mathtt{-r} =$

209 $400$ and $\mathtt{nsites} = 10^6$. From the simulated data sets, we removed all singleton sites to ensure

210 that the sequence data were truly imputable. Within a data set, if the $n = 2000$ haplotypes

211 contained $q$ polymorphic sites after excluding the singletons, we randomly selected, without

212 replacement, $s = 300$ of the $q$ sites, each with minor allele frequency (MAF) greater than

213 0.1. These markers were treated as genotyped. The remaining $q - s$ sites were masked.

²¹⁴ Following ZHANG *et al.* (2013), we calculated the pairwise Hamming distances between the
²¹⁵ $n = 2000$ haplotypes in each of the 50 data sets, based on the genotype information at only
²¹⁶ the $s = 300$ randomly-selected markers. With these distances, we then used the software
²¹⁷ `rapidnj` (SIMONSEN *et al.* 2008) to construct a neighbor-joining tree (SAITOU AND NEI
²¹⁸ 1987) of the haplotypes. Note that it was possible, as a result of random sampling, for two
²¹⁹ or more haplotypes to be identical at all $s$ markers. In such a case, a leaf in the tree would
²²⁰ represent more than one haplotype.

²²¹ Using the python library `dendropy` (SUKUMARAN AND HOLDER 2010), we calculated the
²²² patristic distance matrix for each neighbor-joining tree. We then applied the greedy algo-
²²³ rithm to select the reference panel of size $k = 200$ that maximizes PD. Furthermore, on each
²²⁴ neighbor-joining tree, we used the `rppr` binary in `pplacer` (MATSEN *et al.* 2013) to execute
²²⁵ the adapted PAM algorithm, returning a reference panel of size $k = 200$ that minimizes
²²⁶ ADCL. In cases for which either algorithm selected a leaf that represents more than one
²²⁷ haplotype, one of the haplotypes was randomly chosen to be included in the panel.

²²⁸ In order to model diploid samples, we also created diploid reference panels for use with both
²²⁹ the maximum-PD and minimum-ADCL algorithms. First, we randomly paired the $n = 2000$
²³⁰ haplotypes into 1000 diploid genomes. For the "diploid PD" panel, following ZHANG *et*
²³¹ *al.* (2013), we included diploid individuals carrying at least one of the top-ranked haplotypes
²³² into the panel until we reached the desired panel size $k$. More specifically, we proceeded down
²³³ the list of $k$ haplotypes in the maximum-PD panel, ranked based on the order of inclusion.
²³⁴ At each step, we selected both the top-ranked haplotype and the haplotype with which it
²³⁵ was paired (and which was not necessarily top-ranked) for the diploid panel, if they had not
²³⁶ already been picked previously. We continued this process until $k/2$ diploid genomes were
²³⁷ selected, for a total of $k$ haplotypes.

11

238  Unlike in maximum-PD panels, haplotype sets in minimum-ADCL panels are not nested.

239  Therefore, we cannot use the same process to construct the "diploid ADCL" panel. To

240  address this problem, we first constructed, again using `rppr`, a half-sized minimum-ADCL

241  panel of size $k/2 = 100$. Each haplotype in the half-sized panel, along with the haplotype

242  with which it was paired, was then included in the diploid panel. In the event that both

243  haplotypes of a diploid genome were in the half-sized panel, they were each only chosen once.

244  If the diploid panel was not fully filled at the end of this process, then haplotype pairs were

245  randomly taken from the previously unselected diploid genomes until the requisite panel size

246  of $k/2$ diploid genomes was reached.

247  For comparison, for each of the 50 data sets, we also generated 1000 random reference panels

248  by sampling, without replacement, $k = 200$ of the original $n = 2000$ haplotypes, giving a

249  total of 1004 reference panels. A diagram of the simulation pipeline appears in FIGURE 2.

250  For each of the $k$ haplotypes in a reference panel, we unmasked the genotypes at the $q - s$

251  masked sites and used the resulting full sequences as a reference to perform imputation, under

252  the assumption that the haplotypes represent sequences with resolved phasing. Following

253  ZHANG $et\ al.$ (2013), to avoid edge effects and to improve imputation accuracy, within each

254  1Mb haplotype, we imputed only the middle 100kb segment, while still retaining the markers

255  in both 450kb flanking regions (LI $et\ al.$ 2010). Similar to ZHANG $et\ al.$ (2013), we used

256  the program `minimac` (HOWIE $et\ al.$ 2012) to perform imputation. The parameter values

257  entered into `minimac` were as follows: `--rounds` $= 5$ and `--states` $= 200$.

258  For each choice of reference panel, we evaluated imputation accuracy at the $r$ imputed sites

259  (masked sites within the middle 100kb segment) over the $n/2$ diploid genomes, applying

260  a discordance metric. At imputed site $j$ in diploid genome $i$, we define $g_{ij}$ and $\hat{g}_{ij}$ to be

261  the true and imputed genotypes respectively. Both $g_{ij}$ and $\hat{g}_{ij}$ take on values in $\{0, 1, 2\}$,

12

262 corresponding to the number of copies of an arbitrarily chosen allele at that specific site.

263 The discordance rate $D$ across all sites is given by

$$D = \frac{\sum_{i=1}^{n/2} \sum_{j=1}^{r} |g_{ij} - \hat{g}_{ij}|}{nr}.$$

264 We also compute the discordance rate $H$ across all true heterozygous genotypes ($g_{ij} = 1$):

$$H = \frac{\sum_{i=1}^{n/2} \sum_{j=1}^{r} \mathbf{1}_{g_{ij}=1} |g_{ij} - \hat{g}_{ij}|}{2 \sum_{i=1}^{n/2} \sum_{j=1}^{r} \mathbf{1}_{g_{ij}=1}}.$$

265 In addition, based on the MAF values of their constituent alleles, as computed in the full set

266 of 2000 haplotypes, we further split the true heterozygous sites into three mutually exclusive

267 MAF bins: $0 < \text{MAF} < 0.1$ (low), $0.1 \leq \text{MAF} < 0.2$ (medium), and $0.2 \leq \text{MAF} \leq 0.5$ (high).

268 This separation was performed in order to evaluate how the PD and ADCL algorithms

269 perform across the spectrum of rare to common variants. Note also that the calculations of

270 $D$ and $H$ sum over all $n/2$ diploid genomes, irrespective whether they have one, both, or

271 neither of their haplotypes represented in the reference panel.

272 **1000 Genomes Project sequence data**

273 We also applied both the PD and ADCL algorithms to sequence data from the 1000

274 Genomes Project, available at `http://csg.sph.umich.edu/abecasis/MACH/download/`

275 `1000G-PhaseI-Interim.html`. Following ZHANG *et al.* (2013), we considered $n = 762$

276 phased haplotypes from 381 diploid individuals with European ancestry: 87 Utah residents

277 with Northern and Western European ancestry, 93 Finnish from Finland, 89 British from

278 England and Scotland, 14 Iberians from Spain, and 98 Toscani from Italy.

13

279    We first removed all singleton sites from the data, and we then selected 30 1Mb segments

280    that were approximately evenly spaced across chromosome 20, avoiding the centromere,

281    telomeres, and adjacent areas. Study samples were then created using a similar procedure

282    to that employed for the simulated data. For each of the 30 segments, we randomly selected

283    $s = 400$ markers with MAF $> 0.1$ in the full set of 762 haplotypes, and masked the genotypes

284    of the remaining sites. We then chose $k = 120$ haplotypes to include in the maximum-PD

285    and minimum-ADCL reference panels, as well as in 1000 randomly-generated panels. For

286    each choice of reference panel used for each segment, we imputed the middle 100kb, retaining

287    the markers in both 450kb flanking regions. We then evaluated $D$ and $H$ analogously to the

288    experiments with the simulated data.

## Results

289

### Stability of the adapted PAM algorithm

290

291    Before considering the actual imputation results produced by the different algorithms for

292    reference panel selection, we empirically validated the stability of the adapted PAM algorithm

293    in choosing the minimum-ADCL panel. Beyond the initial run for each of our 50 simulated

294    data sets, we repeated the selection of the minimum-ADCL panel five additional times. For

295    each repetition, we executed the adapted PAM algorithm with a different starting seed, and

296    then determined the number of haplotypes that were shared by the minimum-ADCL panels

297    from both the initial run and the run with the modified seed.

298    When comparing two panels of 200 reference haplotypes drawn from a set of 2000 sample

299    haplotypes, let $m$ be the number of haplotypes that are shared by both panels ($0 \leq m \leq$

300    200). For each of the five replicates, we calculated the mean value of $m$ across the 50

301    data sets, comparing each replicate to the initial run. All five mean values of $m$ were

14

302  observed to be ∼179 (TABLE 1); for comparison, the mean of the hypergeometric distribution

303  describing the number of haplotypes shared between two panels of size 200 independently

304  drawn from a pool of 2000 is 20, with standard deviation 4.03. Therefore, despite changing

305  the specific haplotypes used in randomly initializing the adapted PAM algorithm, most

306  haplotypes eventually chosen for inclusion in the minimum-ADCL panel remain the same.

307  This result suggests that the adapted PAM algorithm is in fact stable, and in subsequent

308  analysis, we consider only a single starting seed.

## Polymorphic sites in reference panels

310  For each of the 1004 reference panels, we evaluated the number of masked sites within the

311  imputed 100kb segment that were polymorphic. This calculation is important because only

312  sites that are polymorphic in the reference panel can produce a meaningful imputation result

313  for the remainder of the study sample. Summing across all 50 data sets, we detected a total

314  of 12,851 masked sites within the 100kb segment of interest. We then compared how many of

315  those masked sites appear as polymorphic in the maximum-PD panel, the minimum-ADCL

316  panel, and a single random panel.

317  Of the 12,851 masked sites, 8879 sites (69.09%) were polymorphic in all three reference-panel

318  types. Of the 3972 remaining sites, 1138 (8.86%) were polymorphic in both the maximum-PD

319  and minimum-ADCL panels, 244 (1.90%) were polymorphic in both the maximum-PD and

320  random panels, and 374 (2.91%) were polymorphic in both the minimum-ADCL and random

321  panels. In addition, 464 (3.61%), 473 (3.68%), and 391 (3.04%) sites were polymorphic in

322  only the maximum-PD, minimum-ADCL, and random panels, respectively. Finally, 888

323  (6.91%) of the masked sites were monomorphic in all three panels (FIGURE 3).

324  Overall, 10,725 sites (83.46%) were polymorphic in the 50 maximum-PD panels, 10,864 sites

325  (84.54%) were polymorphic in the 50 minimum-ADCL panels, and 9888 sites (76.94%) were

326 polymorphic in the 50 random panels. Using the two-tailed Wilcoxon signed-rank test, we

327 found that both the maximum-PD and minimum-ADCL methods of panel selection identify

328 substantially more polymorphic sites compared to choosing the reference panel randomly

329 ($P = 7.686 \times 10^{-10}$ and $P = 8.175 \times 10^{-10}$, respectively).

### Polymorphic sites in imputed data sets

331 The maximum-PD and minimum-ADCL selection algorithms result in similar numbers of

332 polymorphic sites as a fraction of the total number of masked sites in their respective ref-

333 erence panels. We next evaluated the number of imputed sites the two methods recovered

334 as polymorphic. In each of the 50 simulated data sets, we calculated the percentage of

335 masked sites that were polymorphic in the imputed sample, using the maximum-PD panel,

336 the minimum-ADCL panel, the diploid PD panel, the diploid ADCL panel, and the same

337 random panel used to assess the number of polymorphic sites within the reference panels.

338 FIGURE 4 compares the proportion of polymorphic sites imputed with combinations of the

339 five reference panel types. In each panel of FIGURE 4, the random panel is used as a baseline

340 for evaluating two of the other four panel selection methods.

341 We used the two-tailed Wilcoxon signed-rank test to evaluate differences in the fraction

342 of sites identified as polymorphic by the different panel types. Both the maximum-PD

343 and minimum-ADCL panels recover a significantly larger percentage of polymorphic sites

344 compared with their respective diploid panels ($P = 3.448 \times 10^{-9}$ and $P = 2.309 \times 10^{-9}$,

345 respectively). The minimum-ADCL panel also outperforms the maximum-PD panel ($P =$

346 $4.944 \times 10^{-4}$). However, the percentage of imputed sites that are polymorphic shows no

347 significant difference when comparing the diploid PD and diploid ADCL panels ($P = 0.1625$).

16

## Discordance rates

As a measure of imputation accuracy, for each of the 50 simulated data sets, we separately calculated the discordance rate $D$ across all sites that were imputed with the maximum-PD panel, the minimum-ADCL panel, the diploid PD panel, and the diploid ADCL panel. For a baseline, we also calculated the mean discordance rate over the 1000 randomly-selected reference panels. We are mainly interested in comparing the performance between the maximum-PD and minimum-ADCL panels, as well as between the diploid PD and diploid ADCL panels.

The discordance rates appear in FIGURE 5, and their mean values are summarized in TABLE 2. Again using the two-tailed Wilcoxon signed-rank test, the minimum-ADCL panel exhibits significantly lower discordance rates than the maximum-PD panel ($P = 1.342 \times 10^{-9}$). The diploid ADCL panel also has lower discordance rates than the diploid PD panel ($P = 2.597 \times 10^{-3}$). The minimum-ADCL, maximum-PD, diploid ADCL, and diploid PD panels all provide lower discordance rates than the mean of the 1000 randomly-selected panels ($P = 7.789 \times 10^{-10}$, $9.928 \times 10^{-10}$, $8.797 \times 10^{-10}$, and $4.920 \times 10^{-7}$, respectively).

To generate a discordance measure for low-frequency variants, we also calculated the discordance rate $H$ across the heterozygous sites with $0 < \text{MAF} < 0.1$. From FIGURE 5 and TABLE 2, we observe that the mean discordance rates are higher for low-MAF loci than they are for high-MAF loci. Nevertheless, compared to the maximum-PD panel, the minimum-ADCL panel still achieves significantly higher imputation accuracy on low-MAF heterozygotes ($P = 1.606 \times 10^{-9}$). The same relationship also holds between the diploid ADCL and diploid PD panels ($P = 1.871 \times 10^{-4}$). As was observed when considering all variants, the minimum-ADCL, maximum-PD, diploid ADCL, and diploid PD panels all have lower discordance rates than the mean of the 1000 random panels ($P = 7.790 \times 10^{-10}$, $1.264 \times 10^{-9}$, $7.790 \times 10^{-10}$, and $2.244 \times 10^{-6}$, respectively).

17

### Discordance rates under different simulation settings

Following ZHANG *et al.* (2013), to investigate how different parameter choices might have affected the simulation results, we repeated the analysis taking into consideration (i) different reference panel sizes $k$, (ii) different marker densities $s$, and (iii) different target sequence lengths. When varying a parameter, we kept the other two parameters constant at their default values used in the initial analysis (reference panel size $k = 200$, number of markers per MB $s = 300$, imputation length $= 100$kb). The baseline for comparison here is the mean discordance rate over the 50 randomly-selected reference panels. Owing to runtime considerations, this number is smaller than the 1000 randomly-selected reference panels used to calculate the baseline mean discordance rate in the initial analysis. Box plots of the results are shown in FIGURE 6, and mean discordance rates of the various panel types over all sites and over the low-frequency variants appear in TABLES 3 and 4, respectively.

We first evaluated the influence of reference panel size on imputation accuracy, considering cases with $k$ equal to 100, 300, 400, and 500 (compared to the initial analysis with $k = 200$). We observe that as the panel size $k$ increases, discordance rates decrease across all reference panel types. However, we also note a decrease in the difference in performance between the ADCL and PD algorithms, in both the haploid ("maximum-PD" and "minimum-ADCL") and diploid cases. In other words, the gain in imputation accuracy obtained by minimizing ADCL instead of maximizing PD diminishes with large reference panel sizes.

Next, we examined how the initial genotyping density of the markers affected imputation accuracy by considering instances with $s$ equal to 200, 400, 500, and 600 (compared to the initial choice of $s = 300$). Here, across all reference panel types, the discordance rates decrease slightly with increasing marker density $s$. Nevertheless, for all densities, both the haploid and diploid ADCL panels consistently outperform their PD counterparts in terms of imputation accuracy across all sites, as well as across only the low-frequency variants.

18

Finally, we considered whether the length of the target imputation region has an effect on imputation accuracy. We imputed segments of length 500kb, 1Mb and 2Mb (compared to the initial imputation length choice of 100kb). In all cases, a flanking 450kb region was added to each end of the sequence in order to avoid edge effects. We observe that discordance rates remain relatively constant across different imputation lengths. Again, the ADCL panels produce significantly lower discordance rates compared to the PD panels, regardless of the specific choice of imputation length.

## Discordance rates with 1000 Genomes Project sequence data

To confirm that our findings on the simulated data set are also observed when using actual sequence data, we performed a similar analysis for 30 1Mb segments generated on chromosome 20, using 381 diploid individuals with European ancestry from the 1000 Genomes Project. We are again interested in comparing the difference in imputation accuracy achieved by the minimum-ADCL and maximum-PD panels, using the mean discordance rate over 1000 randomly-selected reference panels as a baseline for comparison. The discordance rates appear in FIGURE 7, and their mean values are summarized in TABLE 5. For the three different panel types, FIGURE 8 compares the discordance rates examined in each of the 30 segments over all imputed sites, as well as over only the low-frequency variants.

Applying the two-tailed Wilcoxon signed-rank test, we observe that across all imputed sites, the minimum-ADCL algorithm produces significantly lower discordance rates than the maximum-PD algorithm ($P = 2.367 \times 10^{-3}$), as shown in TABLE 5. In addition, when focusing solely on the low-frequency variants, the minimum-ADCL panel continues to produce better imputation accuracy than the maximum-PD panel ($P = 0.0234$).

19

## Discussion

<sup>419</sup>

<sup>420</sup> The decreasing cost of modern sequencing has enhanced the practicality of generating a

<sup>421</sup> reference panel from the haplotypes that are already present in the study sample. It generally

<sup>422</sup> remains prohibitive, however, to perform full sequencing for large numbers of haplotypes.

<sup>423</sup> Given this constraint in resources, what is the optimal approach for selecting the subset of

<sup>424</sup> the study sample to sequence in order to achieve the best imputation results? We explored

<sup>425</sup> two objective functions for optimization, with the aim of ensuring high imputation accuracy.

<sup>426</sup> Maximizing PD as a way of ensuring that the total genetic diversity of a sample is well-

<sup>427</sup> represented is one sensible approach. This type of panel selection method achieves lower

<sup>428</sup> imputation discordance rates than assembling reference panels from randomly-selected hap-

<sup>429</sup> lotypes (KANG AND MARJORAM 2012; ZHANG *et al.* 2013). Nevertheless, it has not been

<sup>430</sup> clear whether PD represents the best objective function for panel selection.

<sup>431</sup> Minimizing ADCL attempts to ensure that the subset of the study sample selected for the

<sup>432</sup> panel is representative of the total diversity present, albeit using a different approach. It

<sup>433</sup> is conceptually similar to a clustering problem, in that the number of clusters is predeter-

<sup>434</sup> mined, and the algorithm returns the cluster to which each haplotype belongs, as well as the

<sup>435</sup> haplotype that is the most central within its cluster. This haplotype is then included in the

<sup>436</sup> reference panel. Unlike when maximizing PD, the problem of selecting non-representative

<sup>437</sup> branches is mostly avoided by ADCL, as those haplotypes are unlikely to occupy a central

<sup>438</sup> position within their clusters.

<sup>439</sup> For both simulated and actual sequence data, we observed that minimizing ADCL does in fact

<sup>440</sup> provide an improvement in imputation accuracy compared to maximizing PD. It generally

<sup>441</sup> identified a greater number of polymorphic sites, both in the reference panels as well as

<sup>442</sup> in the imputed data. When looking at the overall discordance-rate measures, minimizing

20

443 ADCL produces a significantly lower discordance rate over all sites compared to maximizing

444 PD. This result holds across various choices of genotyping density and imputation length,

445 suggesting that the observed result is robust to such changes. It is only with increasing

446 panel sizes that the gain in imputation accuracy obtained by minimizing ADCL decreases

447 compared to maximizing PD. This outcome could potentially be due to the diminishing

448 returns, in terms of representative variants, contributed by each additional haplotype in the

449 reference panel. Consider the extreme case, where all the haplotypes in the study sample

450 are included in the reference panel. In such a situation, both algorithms return trivially

451 identical imputation results.

452 One metric that is of particular interest is the performance of an algorithm in the imputation

453 of low-frequency variants. Although early genome-wide association (GWA) studies focused

454 on identifying common variants associated with particular diseases or phenotypic traits, the

455 focus of GWA studies has increasingly shifted toward an interest in rare genetic variants

456 (ASIMIT AND ZEGGINI 2010; CIRULLI AND GOLDSTEIN 2010; EICHLER *et al.* 2010). As

457 such studies improve in their ability to detect the effects of rare variants on phenotype (LI *et*

458 *al.* 2013; LEE *et al.* 2014), it is paramount that the imputation process carried out alongside

459 them generate reasonably accurate imputed genotypes with low-frequency variants.

460 In this context, from TABLES 2 and 5, we observed, based on differences in the mean

461 discordance rates, that minimizing ADCL improves upon maximizing PD by the largest

462 absolute amount in the low-MAF bin $(0 < \text{MAF} < 0.1)$, in both the simulated and the actual

463 data. This result might be explained by the fact that the discordance rates obtained when

464 imputing low-frequency variants are relatively high to begin with, and can be potentially

465 reduced to a much greater extent with an improved choice of algorithm for panel selection.

466 Our analyses are consistent in suggesting that the minimum-ADCL algorithm can contribute

467 to reducing imputation inaccuracies in GWA studies that seek to identify the effects of low-

468 frequency variants on phenotypic traits.

21

469 In summary, we have demonstrated that internal reference panel selection via minimizing

470 ADCL produces empirically improved imputation accuracy compared to maximizing PD,

471 particularly for low-frequency variants. This finding applies to both simulated and actual

472 sequence data, and is robust to changes in the choice of initial parameter values. Note

473 that both ADCL and PD represent intermediate criteria that provide practical objective

474 functions, where the ultimate goal is maximizing imputation accuracy or other aspects of

475 imputation performance. Although both algorithms produce considerably better imputation

476 performance measures than the use of random panels, neither is guaranteed to produce the

477 maximal value of such measures over all possible panels. It remains to be determined whether

478 a single simple criterion exists that could lead to identification of the best possible panel for

479 maximizing imputation performance.

## Acknowledgments

## Literature Cited

483 Asimit J., and E. Zeggini, 2010. Rare variant association analysis methods for complex

484 traits. *Annu. Rev. Genet.* **44**: 293–308.

485 Bordewich M., A. G. Rodrigo, and C. Semple, 2008. Selecting a taxa to save or sequence:

486 desirable criteria and a greedy solution. *Syst. Biol.* **57**: 825–834.

487 Cirulli E. T., and D. B. Goldstein, 2010. Uncovering the roles of rare variants in common

488 disease through whole-genome sequencing. *Nat. Rev. Genet.* **11**: 415–425.

Eichler E. E., J. Flint, G. Gibson, A. Kong, S. M. Leal *et al.*, 2010. Missing heritability and strategies for finding the underlying causes of complex disease. *Nat. Rev. Genet.* **11**: 446–450.

Faith D. P., 1992. Conservation evaluation and phylogenetic diversity. *Biol. Conserv.* **61**: 1–10.

Fridley B. L., G. Jenkins, M. E. Deyo-Svendsen, S. Hebbring, and R. Freimuth, 2010. Utilizing genotype imputation for the augmentation of sequence data. *PLoS ONE* **5**: e11018.

Hartmann K., and M. Steel, 2007. Phylogenetic diversity: from combinatorics to ecology, pp. 171–196 in *Reconstructing Evolution: New Mathematical and Computational Advances*, edited by O. Gascuel, and M. Steel. Oxford University Press, Oxford.

Howie, B., C. Fuchsberger, M. Stephens, J. Marchini, and G. R. Abecasis, 2012. Fast and accurate genotype imputation in genome-wide association studies through pre-phasing. *Nat. Genet.* **44**: 955–959.

Huang, L., Y. Li, A. B. Singleton, J. A. Hardy, G. Abecasis *et al.*, 2009. Genotype-imputation accuracy across worldwide human populations. *Am. J. Hum. Genet.* **84**: 235–250.

Huang, L., M. Jakobsson, T. J. Pemberton, M. Ibrahim, T. Nyambo *et al.*, 2011. Haplotype variation and genotype imputation in African populations. *Genet. Epidemiol.* **35**: 766–780.

Hudson, R. R., 2002. Generating samples under a Wright-Fisher neutral model of genetic variation. *Bioinformatics* **18**: 337–338.

International HapMap Consortium, 2005. A haplotype map of the human genome. *Nature* **437**: 1299–1320.

Jewett, E. M., M. Zawistowski, N. A. Rosenberg, and S. Zöllner, 2012. A coalescent model for genotype imputation. *Genetics* **191**: 1239–1255.

Kang, C. J., and P. Marjoram, 2012. A sample selection strategy for next-generation sequencing. *Genet. Epidemiol.* **36**: 696–709.

Kaufman, L., and P. J. Rousseeuw, 1987. Clustering by means of medoids, pp. 405–416 in *Statistical Data Analysis Based on the L1-Norm and Related Methods*, edited by Y. Dodge. North-Holland, Amsterdam.

Kreiner-Møller, E., C. Medina-Gomez, A. G. Uitterlinden, F. Rivadeneira, and K. Estrada, 2015. Improving accuracy of rare variant imputation with a two-step imputation approach. *Eur. J. Hum. Genet.* **23**: 395–400.

Lee, S., G. R. Abecasis, M. Boehnke, and X. Lin, 2014. Rare-variant association analysis: study designs and statistical tests. *Am. J. Hum. Genet.* **95**: 5–23.

Li, B., D. J. Liu, and S. M. Leal, 2013. Identifying rare variants associated with complex traits via sequencing. *Curr. Protoc. Hum. Genet.* **78**: 1.26.1–1.26.22.

Li, Y., C. Willer, S. Sanna, and G. Abecasis, 2009. Genotype imputation. *Annu. Rev. Genomics Hum. Genet.* **10**: 387–406.

Li, Y., C. J. Willer, J. Ding, P. Scheet, and G. R. Abecasis, 2010. MaCH: using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genet. Epidemiol.* **34**: 816–834.

Marchini, J., and B. Howie, 2010. Genotype imputation for genome-wide association studies. *Nat. Rev. Genet.* **11**: 499–511.

Matsen, F. A., A. Gallagher, and C. O. McCoy, 2013. Minimizing the average distance to a closest leaf in a phylogenetic tree. *Syst. Biol.* **62**: 824–836.

Pardi, F., and N. Goldman, 2005. Species choice for comparative genomics: being greedy works. *PLoS Genet.* **1**: e71.

Paşaniuc, B., R. Avinery, T. Gur, C. F. Skibola, P. M. Bracci *et al.*, 2010. A generic coalescent-based framework for the selection of a reference panel for imputation. *Genet. Epidemiol.* **34**: 773–782.

Peil, B., M. Kabisch, C. Fischer, U. Hamann, and J. L. Bermejo, 2015. Tailored selection of study individuals to be sequenced in order to improve the accuracy of genotype imputation. *Genet. Epidemiol.* **39**: 114–121.

Saitou, N., and M. Nei, 1987. The neighbor-joining method: a new method for reconstructing phylogenetic tress. *Mol. Biol. Evol.* **4**: 406–425.

Sampson J. N., K. Jacobs, Z. Wang, M. Yeager, S. Chanock *et al.*, 2012. A two-platform design for next generation genome-wide association studies. *Genet. Epidemiol.* **36**: 400–408.

Sheng, W., and X. Liu, 2004. A hybrid algorithm for $k$-medoid clustering of large data sets. *Proc. IEEE Congr. Evol. Comput.* **1**: 77–82.

Shriner, D., and A. Adeyemo, G. Chen, C. N. Rotimi, 2010. Practical considerations for imputation of untyped markers in admixed populations. *Genet. Epidemiol.* **34**: 258–265.

Simonsen, M., T. Mailund, and C. N. S. Pedersen, 2008. Rapid neighbor-joining, pp. 113–122 in *Algorithms in Bioinformatics*, edited by K. A. Crandall, and J. Lagergren. Springer-Verlag, Berlin.

Steel, M., 2005. Phylogenetic diversity and the greedy algorithm. *Syst. Biol.* **54**: 527–529.

553  Sukumaran, J., and M. T. Holder, 2010. DendroPy: A Python library for phylogenetic

554  computing. *Bioinformatics* **26**: 1569–1571.

555  Surakka, I., K. Kristiansson, V. Anttila, M. Inouye, C. Barnes *et al.*, 2010. Founder

556  population-specific HapMap panel increases power in GWA studies through improved impu-

557  tation accuracy and CNV tagging. *Genome Res.* **20**: 1344–1351.

558  The 1000 Genomes Project Consortium, 2010. A map of human genome variation from

559  population-scale sequencing. *Nature* **467**: 1061–1073.

560  Theodoridis, S., and K. Koutroumbas, 2008. *Pattern Recognition* (4th ed.). Academic Press,

561  Waltham, MA.

562  Zhang, P., X. Zhan, N. A. Rosenberg, and S. Zöllner, 2013. Genotype imputation reference

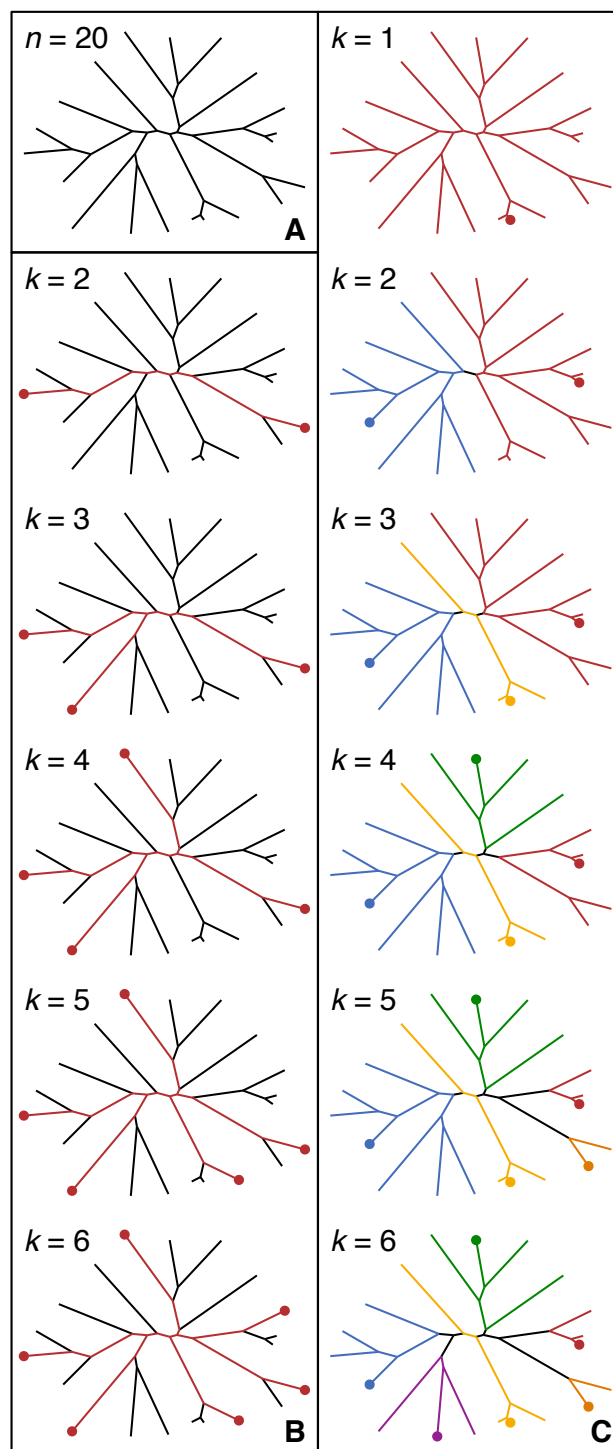563  panel selection using maximal phylogenetic diversity. *Genetics* **195**: 319–330.

FIGURE 1: Reference panels for an example tree with $n = 20$ haplotypes. (A) An example tree. (B) The maximum-PD panel. (C) The minimum-ADCL panel. In (B) and (C), the haplotypes selected for a given panel size $k$ are represented by a dot at the tips. In (C), each selected haplotype is assigned a color, and all other branches share a color with the closest selected haplotype.
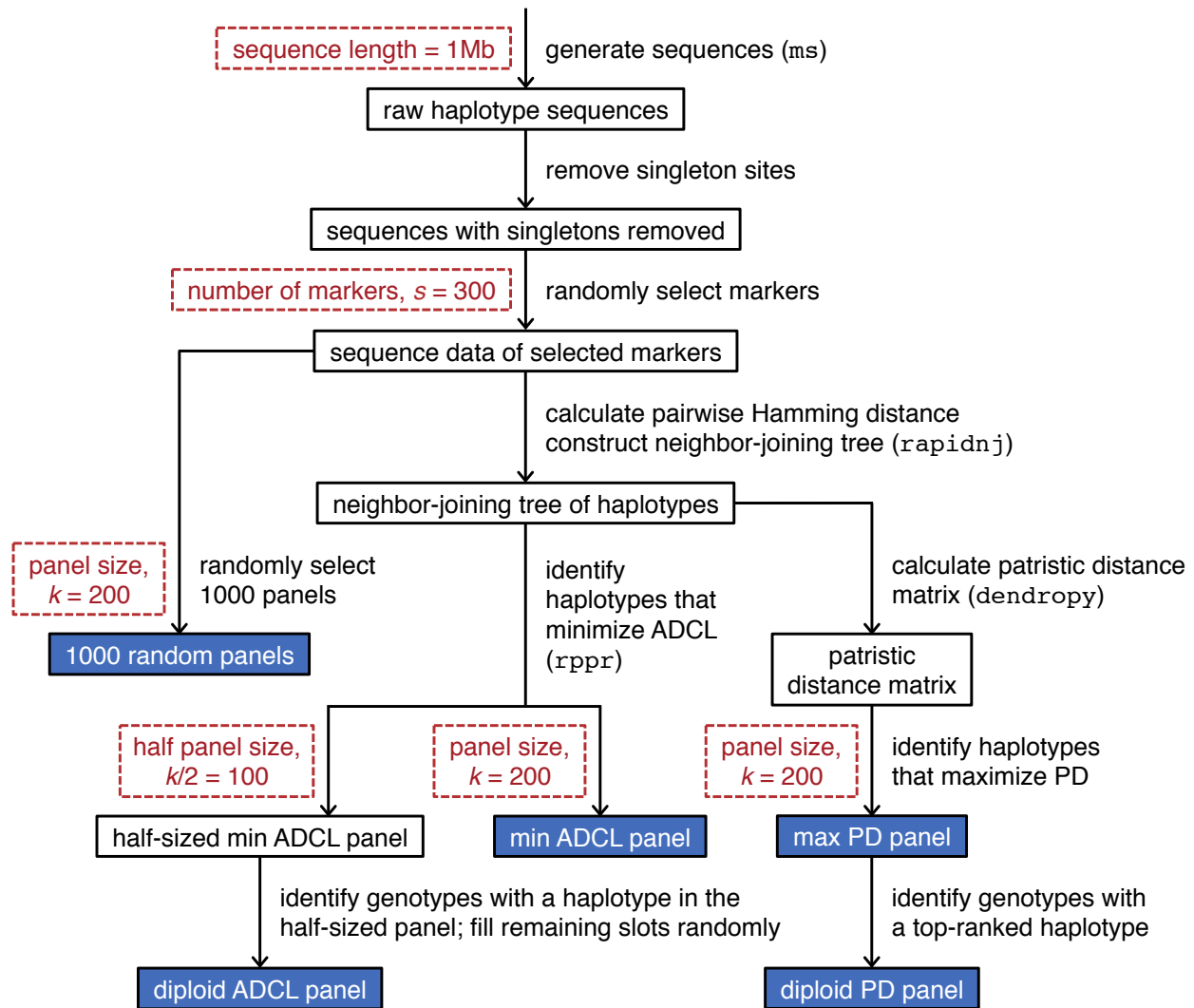
FIGURE 2: A schematic diagram of the pipeline used to generate the simulated data. The red boxes each represent a parameter choice, and the blue boxes represent the 1004 reference panels used in our evaluation.
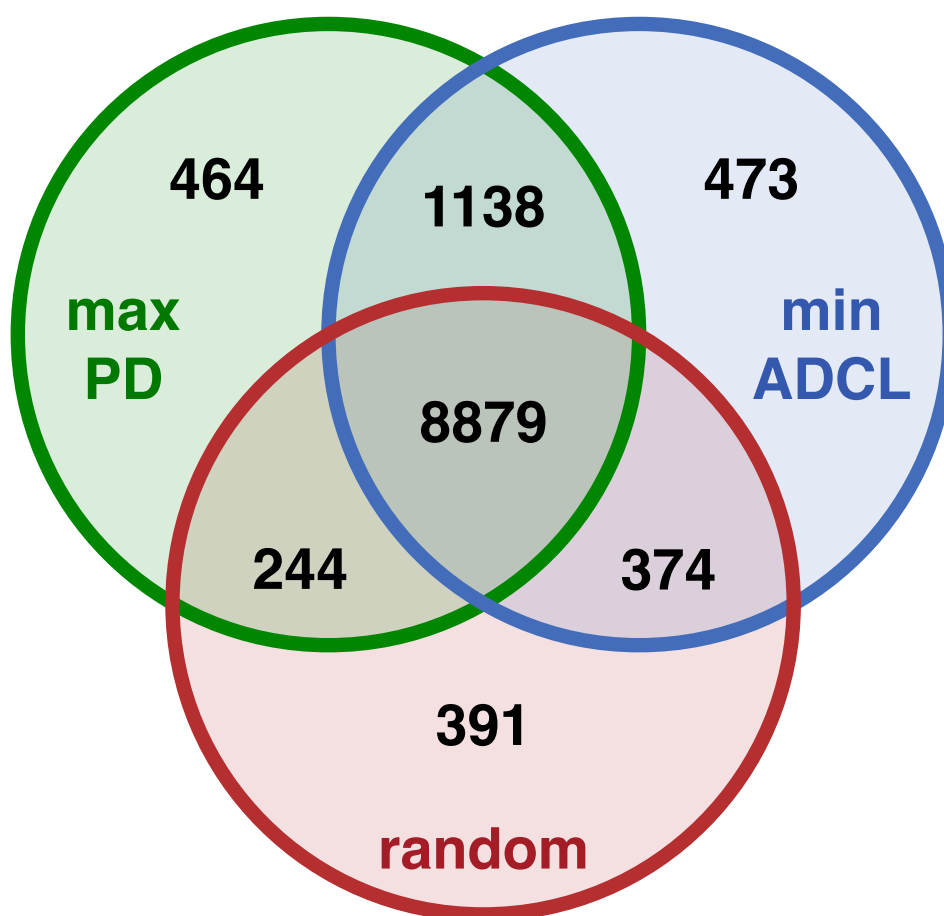
FIGURE 3: A Venn diagram showing the number of polymorphic sites returned by each panel type, out of a total of 12,851 masked sites. 888 sites were monomorphic in all three panels.
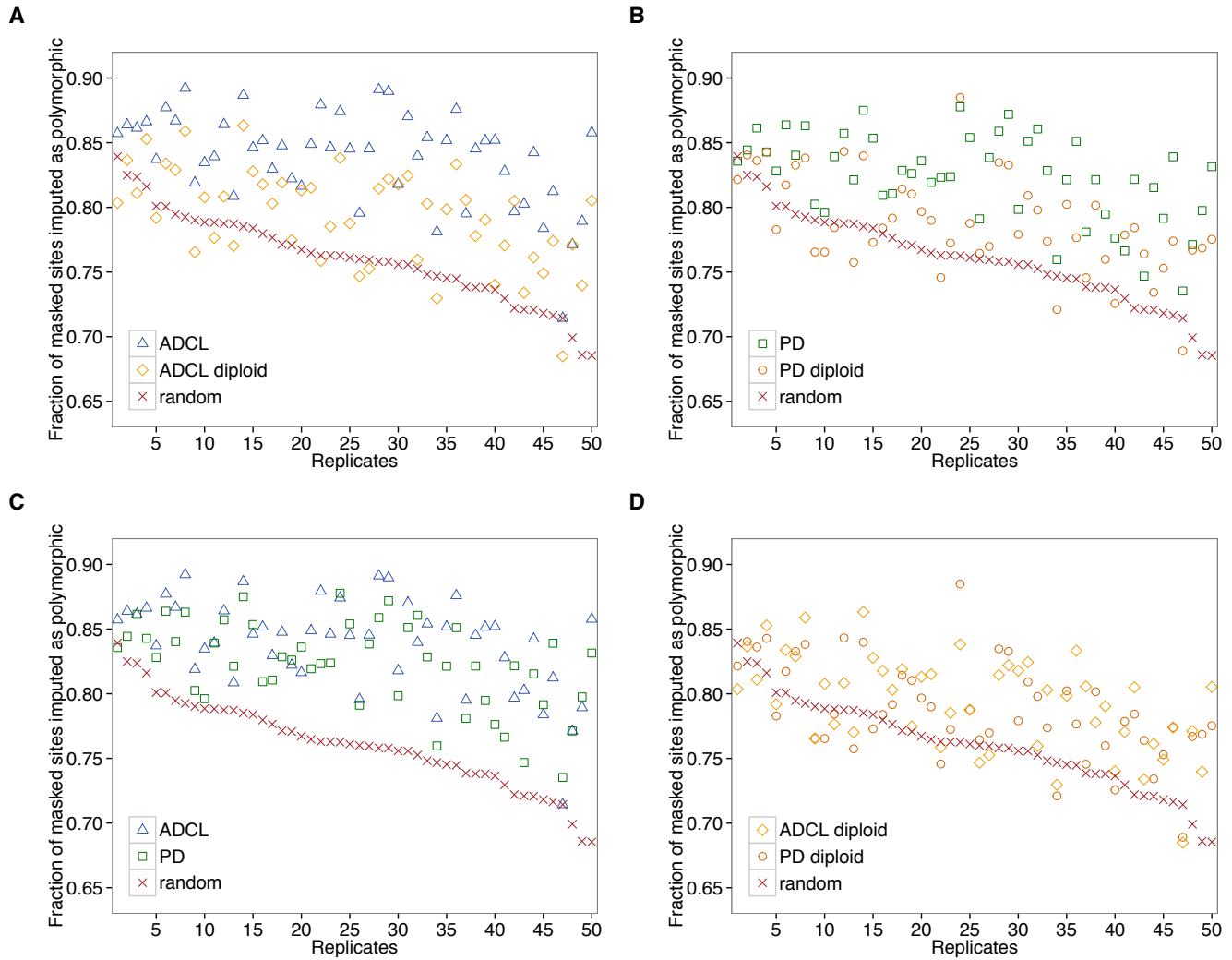
FIGURE 4: Fraction of masked sites imputed as polymorphic, using five different types of reference panels. Data are split into various graphs for ease of comparison. (**A**) ADCL versus ADCL diploid. (**B**) PD versus PD diploid. (**C**) ADCL versus PD. (**D**) ADCL diploid versus PD diploid. The 50 replicate data sets are sorted in decreasing order by the percentage of polymorphic sites recovered by imputations using the random reference panel.
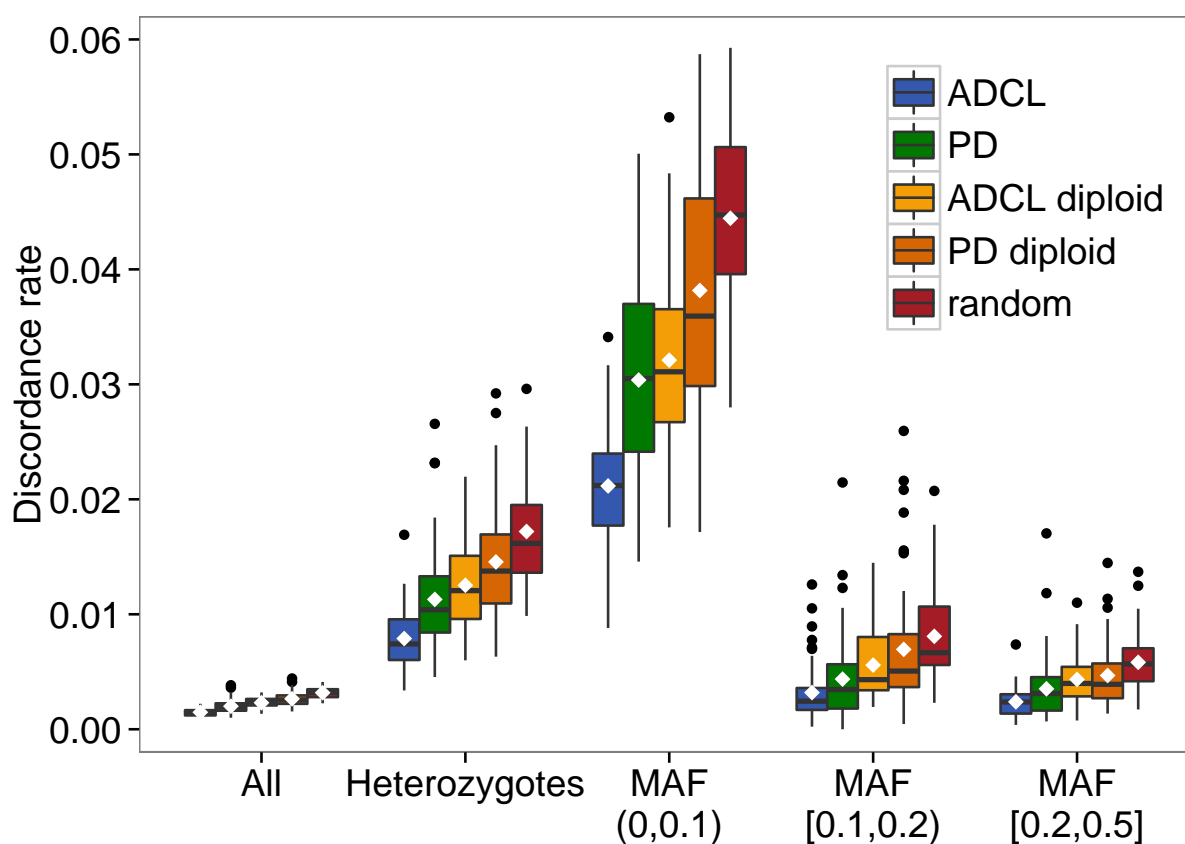
FIGURE 5: Box plots of discordance rates between imputed and simulated genotypes using the five different reference panel types. The mean discordance rate across the 50 replicates for each comparison group is indicated by a diamond, and the median discordance rate is indicated by a horizontal line. The x-axis separates the comparison over all sites, all heterozygous sites, and heterozygous sites falling into three different MAF groups.
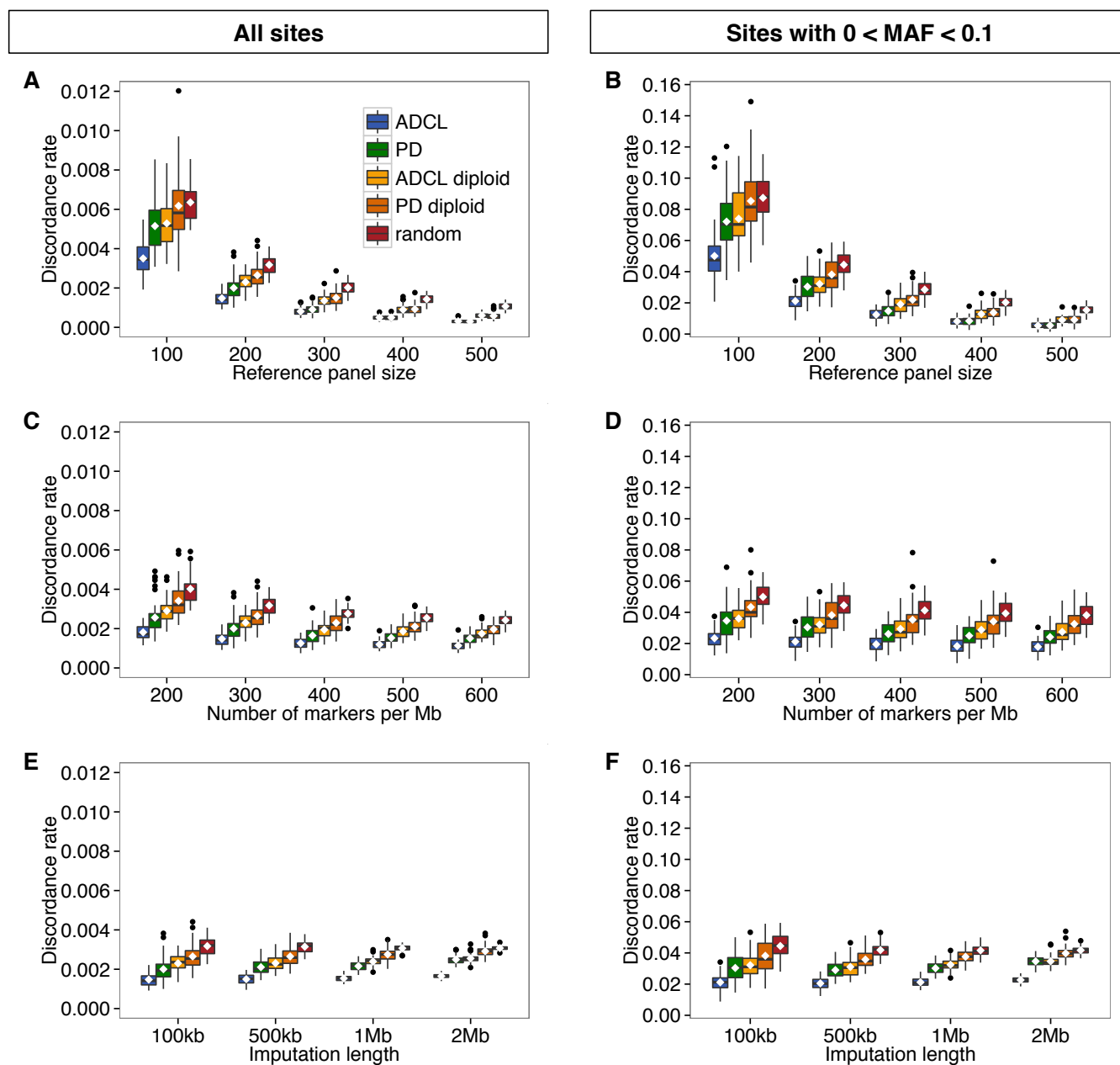
FIGURE 6: Box plots of discordance rates between imputed and simulated genotypes using the five different reference panel types. (**A**) Varying reference panel size, all sites. (**B**) Varying reference panel size, heterozygous sites with $0 < \mathrm{MAF} < 0.1$. (**C**) Varying marker density, all sites. (**D**) Varying marker density, heterozygous sites with $0 < \mathrm{MAF} < 0.1$. (**E**) Varying imputation length, all sites. (**F**) Varying imputation length, heterozygous sites with $0 < \mathrm{MAF} < 0.1$.
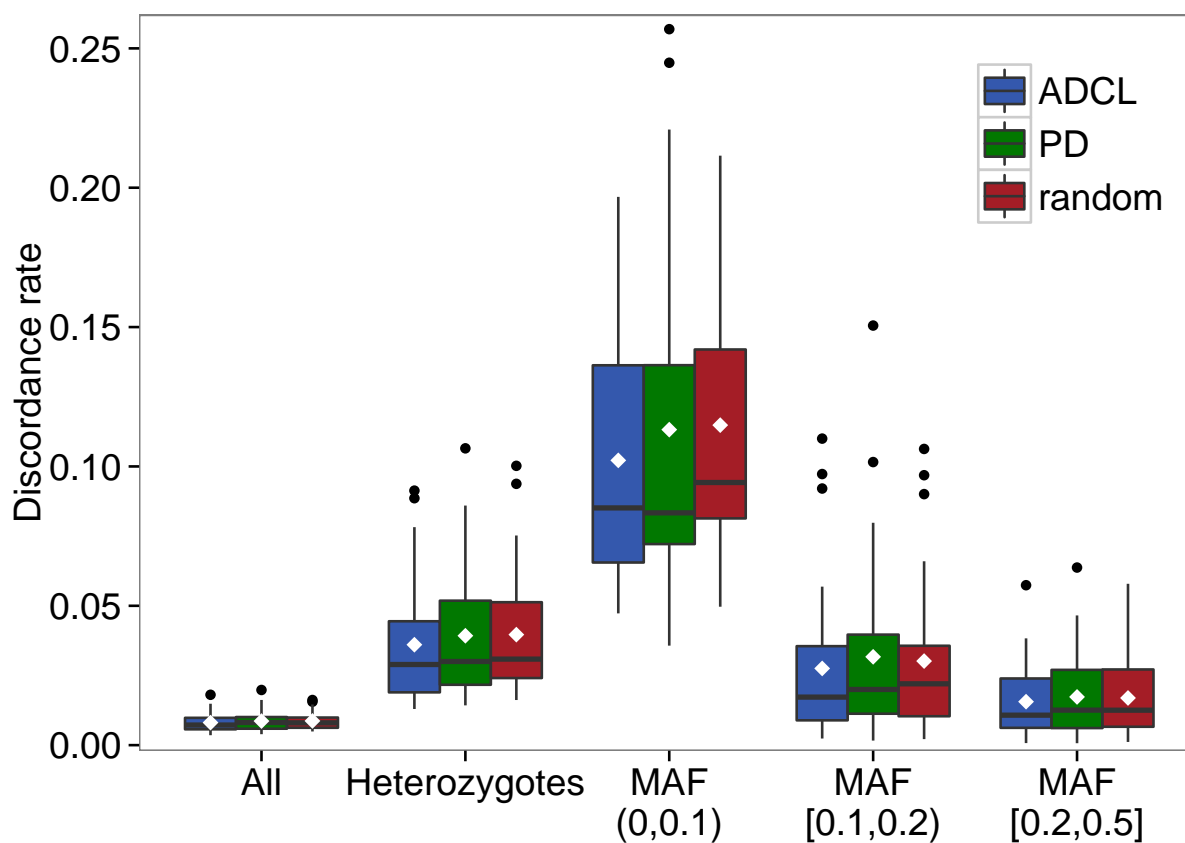
FIGURE 7: Box plots of discordance rates between imputed and actual genotypes using the minimum-ADCL, maximum-PD, and random panels. The data consist of 30 1Mb segments from 762 haplotypes of European ancestry obtained from the 1000 Genomes Project. The mean discordance rate across the 30 replicates for each comparison group is indicated by a diamond, and the median discordance rate is indicated by a horizontal line. The $x$-axis separates the comparison over all sites, all heterozygous sites, and heterozygous sites falling into three different MAF groups.
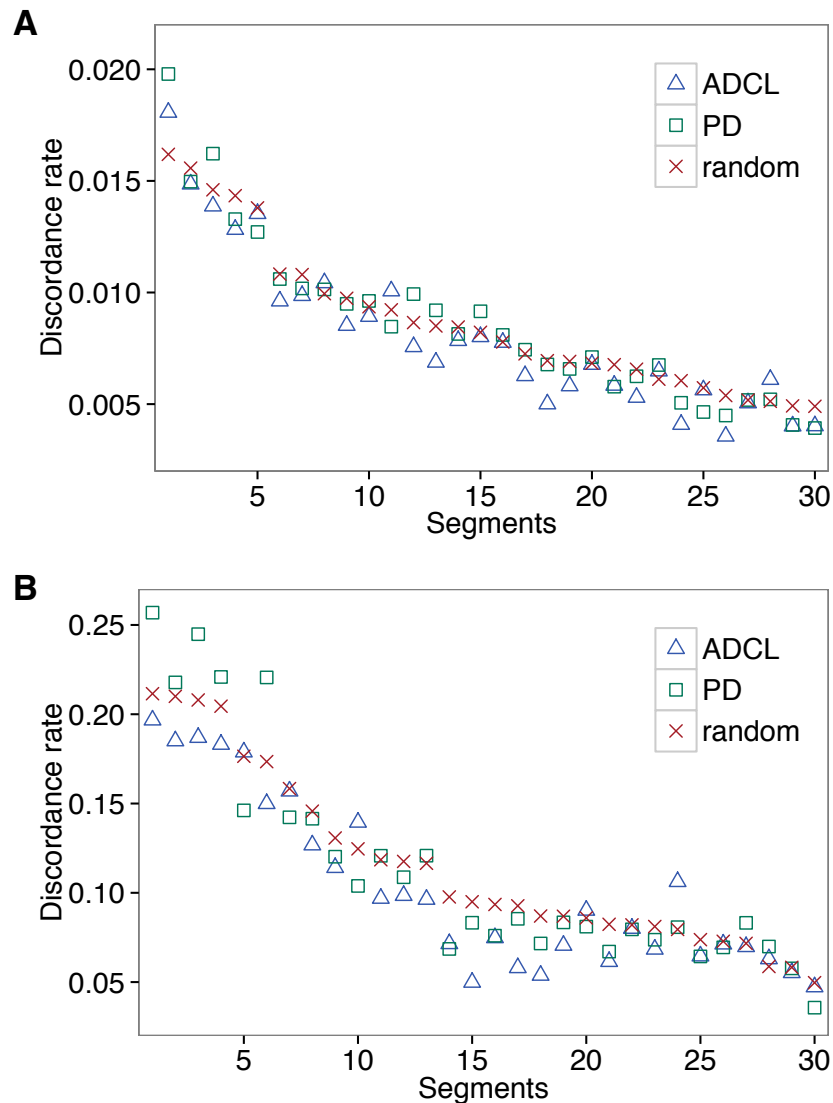
FIGURE 8: Discordance rates between imputed and actual genotypes using the minimum-ADCL, maximum-PD, and random panels, showing an alternative presentation of the same data used to generate FIGURE 7. (A) All sites. (B) Heterozygous sites with $0 < \text{MAF} < 0.1$. The 30 segments are sorted in decreasing order by the mean discordance rate over 1000 random panels.

TABLE 1: Mean and standard deviation of the number of shared haplotypes across 50 data sets in each of five replicates

| Replicate | Mean | Standard deviation |
|:---:|:---:|:---:|
| 1 | 179.40 | 4.1991 |
| 2 | 178.56 | 4.9494 |
| 3 | 179.38 | 4.5888 |
| 4 | 179.00 | 4.7208 |
| 5 | 178.58 | 4.8910 |

For the five replicates, each with a different starting seed, we compared the minimum-ADCL panels from the initial run of the adapted PAM algorithm and the minimum-ADCL panels using the different seed. The table shows the mean (out of 200) and standard deviation of the number of shared haplotypes across the 50 data sets.

TABLE 2: Mean discordance rates between imputed and simulated genotypes, using the maximum (haploid) PD, minimum (haploid) ADCL, diploid PD, and diploid ADCL panels

| | Haploid panels | | | Diploid panels | | |
|---|---|---|---|---|---|---|
| | PD (%) | ADCL (%) | $P$-value | PD (%) | ADCL (%) | $P$-value |
| All | 0.2003 | 0.1476 | $1.342 \times 10^{-9}$ | 0.2648 | 0.2304 | $2.597 \times 10^{-3}$ |
| Heterozygotes | 1.1295 | 0.7888 | $1.921 \times 10^{-9}$ | 1.4554 | 1.2523 | $1.360 \times 10^{-3}$ |
| MAF $(0, 0.1)$ | 3.0374 | 2.1160 | $1.606 \times 10^{-9}$ | 3.8164 | 3.2103 | $1.871 \times 10^{-4}$ |
| MAF $[0.1, 0.2)$ | 0.4363 | 0.3190 | $2.792 \times 10^{-3}$ | 0.6947 | 0.5582 | 0.0857 |
| MAF $[0.2, 0.5]$ | 0.3518 | 0.2386 | $2.567 \times 10^{-5}$ | 0.4676 | 0.4335 | 0.4174 |

This table is obtained from the data in FIGURE 5. The comparison is performed over all sites, all heterozygous sites, and heterozygous sites falling into three different MAF groups. Also shown are the $P$-values of the two-tailed Wilcoxon signed-rank tests comparing the discordance rates of the PD and ADCL reference panels.

36

TABLE 3: Mean discordance rates between imputed and simulated genotypes for all sites, using the maximum (haploid) PD, minimum (haploid) ADCL, diploid PD, and diploid ADCL panels, under different input parameter choices

| | Haploid panels | | | Diploid panels | | |
|---|---|---|---|---|---|---|
| | PD (%) | ADCL (%) | $P$-value | PD (%) | ADCL (%) | $P$-value |
| Reference panel size, $k$ | | | | | | |
| $k = 100$ | 0.5150 | 0.3502 | $5.213 \times 10^{-9}$ | 0.6174 | 0.5284 | $2.257 \times 10^{-5}$ |
| $\boldsymbol{k = 200}$ | **0.2003** | **0.1476** | $\mathbf{1.342 \times 10^{-9}}$ | **0.2648** | **0.2304** | $\mathbf{2.597 \times 10^{-3}}$ |
| $k = 300$ | 0.0907 | 0.0811 | $2.767 \times 10^{-3}$ | 0.1501 | 0.1354 | 0.0167 |
| $k = 400$ | 0.0499 | 0.0498 | 0.7391 | 0.0924 | 0.0895 | 0.3417 |
| $k = 500$ | 0.0298 | 0.0312 | 0.2112 | 0.0584 | 0.0605 | 0.1765 |
| Number of markers per MB, $s$ | | | | | | |
| $s = 200$ | 0.2561 | 0.1810 | $1.378 \times 10^{-8}$ | 0.3409 | 0.2901 | $1.173 \times 10^{-4}$ |
| $\boldsymbol{s = 300}$ | **0.2003** | **0.1476** | $\mathbf{1.342 \times 10^{-9}}$ | **0.2648** | **0.2304** | $\mathbf{2.597 \times 10^{-3}}$ |
| $s = 400$ | 0.1647 | 0.1255 | $2.548 \times 10^{-8}$ | 0.2291 | 0.1946 | $4.176 \times 10^{-5}$ |
| $s = 500$ | 0.1529 | 0.1193 | $9.347 \times 10^{-10}$ | 0.2105 | 0.1863 | $7.284 \times 10^{-4}$ |
| $s = 600$ | 0.1503 | 0.1138 | $4.130 \times 10^{-9}$ | 0.1970 | 0.1746 | $6.975 \times 10^{-5}$ |
| Imputation length | | | | | | |
| **100kb** | **0.2003** | **0.1476** | $\mathbf{1.342 \times 10^{-9}}$ | **0.2648** | **0.2304** | $\mathbf{2.597 \times 10^{-3}}$ |
| 500kb | 0.2104 | 0.1494 | $7.789 \times 10^{-10}$ | 0.2650 | 0.2307 | $5.575 \times 10^{-6}$ |
| 1Mb | 0.2159 | 0.1538 | $7.790 \times 10^{-10}$ | 0.2755 | 0.2392 | $2.040 \times 10^{-8}$ |
| 2Mb | 0.2495 | 0.1653 | $7.789 \times 10^{-10}$ | 0.2914 | 0.2551 | $8.263 \times 10^{-9}$ |

The table is obtained from the data in FIGURES 6A, C and E. Also shown are the $P$-values of the two-tailed Wilcoxon signed-rank tests comparing the discordance rates of the PD and ADCL reference panels. The discordance rates and $P$-values from the initial analysis using $k = 200$, $s = 300$ and imputation length = 100kb are given in bold, with the values obtained from TABLE 2.

TABLE 4: Mean discordance rates between imputed and simulated genotypes for all heterozygous sites with $0 < \text{MAF} < 0.1$, using the maximum (haploid) PD, minimum (haploid) ADCL, diploid PD, and diploid ADCL panels, under different input parameter choices

| | Haploid panels | | | Diploid panels | | |
|---|---|---|---|---|---|---|
| | PD (%) | ADCL (%) | $P$-value | PD (%) | ADCL (%) | $P$-value |
| Reference panel size, $k$ | | | | | | |
| $k = 100$ | 7.2099 | 5.0056 | $3.358 \times 10^{-8}$ | 8.5331 | 7.3942 | $3.960 \times 10^{-4}$ |
| $\boldsymbol{k = 200}$ | **3.0374** | **2.1160** | $\mathbf{1.606 \times 10^{-9}}$ | **3.8164** | **3.2103** | $\mathbf{1.871 \times 10^{-4}}$ |
| $k = 300$ | 1.4783 | 1.2425 | $1.484 \times 10^{-4}$ | 2.1964 | 1.9086 | $3.817 \times 10^{-4}$ |
| $k = 400$ | 0.8394 | 0.8156 | 0.2972 | 1.3786 | 1.2816 | 0.0757 |
| $k = 500$ | 0.5494 | 0.5413 | 0.7502 | 0.9160 | 0.9010 | 0.7138 |
| Number of markers per MB, $s$ | | | | | | |
| $s = 200$ | 3.4597 | 2.3336 | $3.467 \times 10^{-9}$ | 4.3339 | 3.5993 | $8.581 \times 10^{-6}$ |
| $\boldsymbol{s = 300}$ | **3.0374** | **2.1160** | $\mathbf{1.606 \times 10^{-9}}$ | **3.8164** | **3.2103** | $\mathbf{1.871 \times 10^{-4}}$ |
| $s = 400$ | 2.6079 | 1.9485 | $3.270 \times 10^{-9}$ | 3.5185 | 2.9342 | $1.173 \times 10^{-5}$ |
| $s = 500$ | 2.4951 | 1.8300 | $1.810 \times 10^{-9}$ | 3.4146 | 2.8719 | $7.874 \times 10^{-5}$ |
| $s = 600$ | 2.4121 | 1.7891 | $7.368 \times 10^{-9}$ | 3.2507 | 2.7438 | $1.528 \times 10^{-5}$ |
| Imputation length | | | | | | |
| **100kb** | **3.0374** | **2.1160** | $\mathbf{1.606 \times 10^{-9}}$ | **3.8164** | **3.2103** | $\mathbf{1.871 \times 10^{-4}}$ |
| 500kb | 2.8998 | 2.0484 | $7.790 \times 10^{-10}$ | 3.5755 | 3.0995 | $1.605 \times 10^{-6}$ |
| 1Mb | 3.0180 | 2.1204 | $7.790 \times 10^{-10}$ | 3.7531 | 3.2439 | $1.231 \times 10^{-8}$ |
| 2Mb | 3.4652 | 2.2648 | $7.790 \times 10^{-10}$ | 3.9769 | 3.4637 | $9.806 \times 10^{-9}$ |

The table is obtained from the data in FIGURES 6B, D and F. Also shown are the $P$-values of the two-tailed Wilcoxon signed-rank tests comparing the discordance rates of the PD and ADCL reference panels. The discordance rates and $P$-values from the initial analysis using $k = 200$, $s = 300$ and imputation length = 100kb are given in bold, with the values obtained from TABLE 2.

TABLE 5: Mean discordance rates between imputed and 1000 Genomes genotypes, using the maximum-PD and minimum-ADCL panels

| | PD (%) | ADCL (%) | $P$-value |
|---|---|---|---|
| All | 0.8643 | 0.8087 | $2.367 \times 10^{-3}$ |
| Heterozygotes | 3.9253 | 3.6041 | $9.301 \times 10^{-3}$ |
| MAF $(0, 0.1)$ | 11.3202 | 10.2176 | 0.0234 |
| MAF $[0.1, 0.2)$ | 3.1695 | 2.7525 | 0.0274 |
| MAF $[0.2, 0.5]$ | 1.7302 | 1.5598 | $8.035 \times 10^{-4}$ |

The table is obtained from the data in FIGURES 7 and 8. The comparison is performed over all sites, all heterozygous sites, and heterozygous sites falling into three different MAF groups. Also shown are the $P$-values of the two-tailed Wilcoxon signed-rank tests comparing the discordance rates of the PD and ADCL reference panels.