

1 Rare variants imputation in admixed populations:
2 Comparison across reference panels and bioinformatics tools.

3
4 Sanjeev Sariya^{a,b}, Joseph Lee^{a,b}, Richard Mayeux^{a,b,c}, Badri N. Vardarajan^{a,b}, Dolly Reyes-
5 Dumeyer^{a,b}, Jennifer Manly, Adam Brickman, Rafael Lantigua^d, Martin Medrano^e, Ivonne Z.
6 Jimenez-Velazquez^f and Giuseppe Tosto^{a,b,c}

- 7
8 a. Taub Institute for Research on Alzheimer's Disease and the Aging Brain, College of
9 Physicians and Surgeons, Columbia University. 630 West 168th Street, New York, NY
10 10032.
11 b. The Gertrude H. Sergievsky Center, College of Physicians and Surgeons, Columbia
12 University. 630 West 168th Street, New York, NY 10032.
13 c. Department of Neurology, College of Physicians and Surgeons, Columbia University and
14 the New York Presbyterian Hospital. 710 West 168th Street, New York, NY 10032
15 d. Medicine College of Physicians and Surgeons, and The Department of 6Epidemiology,
16 School of Public Health, Columbia University, New York, NY, USA
17 e. School of Medicine, Pontificia Universidad Catolica Madre y Maestra, Santiago,
18 Dominican Republic
19 f. Department of Medicine, Geriatrics Program, School of Medicine, University of Puerto
20 Rico, San Juan, Puerto Rico

21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37

Word counts: Abstract 344; Main Text 3,821; Tables: 4; Figures: 2; References: 35

Correspondence:
Giuseppe Tosto MD PhD
Taub Institute for Research on
Alzheimer's Disease and the Aging Brain
630 West 168th Street
New York, NY 10032
Tel: +1 212 305 9274
Email: gt2260@cumc.columbia.edu

38 **Abstract.**

39 **Background:** Imputation has become a standard approach in genome-wide association studies
40 (GWAS) to infer *in silico* untyped markers. Although feasibility for common variants imputation
41 is well established, we aimed to assess rare and ultra-rare variants' imputation in an admixed
42 Caribbean Hispanic population (CH).

43

44 **Methods:** We evaluated imputation accuracy in CH (N=1,000), focusing on rare ($0.1\% \leq$ minor
45 allele frequency (MAF) $\leq 1\%$) and ultra-rare (MAF $< 0.1\%$) variants. We used two reference
46 panels, the Haplotype Reference Consortium (HRC; N=27,165) and 1000 Genome Project
47 (1000G phase 3; N=2,504) and multiple phasing (SHAPEIT, Eagle2) and imputation algorithms
48 (IMPUTE2, MACH-Admix). To assess imputation quality, we reported: a) high-quality variant
49 counts according to imputation tools' internal indexes (e.g. IMPUTE2 "Info" $\geq 80\%$). b)
50 Wilcoxon Signed-Rank Test comparing imputation quality for genotyped variants that were
51 masked and imputed; c) Cohen's kappa coefficient to test agreement between imputed and
52 whole-exome sequencing (WES) variants; d) imputation of G206A mutation in the *PSEN1*
53 (ultra-rare in the general population and more frequent in CH) followed by confirmation
54 genotyping. We also tested ancestry proportion (European, African and Native American)
55 against WES-imputation mismatches in a Poisson regression fashion.

56

57 **Results:** SHAPEIT2 retrieved higher percentage of imputed high-quality variants than Eagle2
58 (rare: 51.02% vs. 48.60%; ultra-rare 0.66% vs 0.65%, Wilcoxon p-value < 0.001). SHAPEIT-
59 IMPUTE2 employing HRC outperformed 1000G (64.50% vs. 59.17%; 1.69% vs 0.75% for high-
60 quality rare and ultra-rare variants, respectively; Wilcoxon p-value < 0.001). SHAPEIT-
61 IMPUTE2 outperformed MaCH-Admix. Compared to 1000G, HRC-imputation retrieved a

62 higher number of high-quality rare and ultra-rare variants, despite showing lower agreement
63 between imputed and WES variants (e.g. rare: 98.86% for HRC vs. 99.02% for 1000G). High
64 Kappa ($K = 0.99$) was observed for both reference panels. Twelve G206A mutation carriers were
65 imputed and all validated by confirmation genotyping. African ancestry was associated with
66 higher imputation errors for uncommon and rare variants (p-value < 1e-05).

67
68 **Conclusion:** Reference panels with larger numbers of haplotypes can improve imputation quality
69 for rare and ultra-rare variants in admixed populations such as CH. Ethnic composition is an
70 important predictor of imputation accuracy, with higher African ancestry associated with poorer
71 imputation accuracy.

72

73

74 **Keywords:** Rare variants, Imputation, Admixed population, GWAS, 1000G

75 **Introduction**

76 Genome-wide association studies (GWASs) are a major tool to identify common variants
77 associated with complex diseases. GWAS can include 550K to over 2M Single Nucleotide
78 Polymorphisms (SNPs) (Ha et al., 2014) to cover the human genome evenly. Although GWAS
79 has shown to be a robust method to identify disease loci of interest, they rarely point to a causal
80 coding variant. In fact, microarray SNP chips for GWAS are optimally designed to uncover
81 common variants, often associated with small effect sizes mostly located in intronic and
82 intergenic regions. The focus of genetic investigations has since shifted toward rarer alleles with
83 larger effect sizes (Gibson, 2012). With the changing paradigm, imputation of rare variants has
84 become an important topic to enhance the genome coverage in GWAS. Imputation is a process
85 of inferring untyped SNP markers in the discovery population by using densely typed SNPs in
86 external reference panel(s). These '*in silico*' markers increase the coverage of association tests
87 while conducting genome-wide association analysis. In addition, large number of SNPs facilitate
88 meta-analysis when merging data from different study cohorts.

89 The quality of imputation essentially depends on two parameters: available reference datasets
90 and algorithms that employ those reference datasets. Previous studies have shown that
91 imputation quality depends on how well reference panels reflect the study population. To
92 respond to the needs, the 1000 Genome project (1000G), now in its third phase release, has
93 proven to be one of the most frequently used reference panels (Genomes Project et al., 2015).
94 Using these composite reference panels, a number of studies (Pei et al., 2010; Howie et al., 2012;
95 Verma et al., 2014; Liu et al., 2015) have compared imputation accuracy using different
96 imputation tools and algorithms, although the results are equivocal. Few studies (Browning and
97 Browning, 2009; Zheng et al., 2012; Zheng et al., 2015) assessed the impact of reference panel

98 size and input data's features - such as density of SNPs - to impute rare variants, suggesting
99 larger size of reference panels work better. Surakka and colleagues (Surakka et al., 2016)
100 assessed accuracy of imputed SNPs by evaluating rate of false polymorphisms in a Finnish
101 population using global reference panels – Haplotype Reference Consortium (HRC) release 1,
102 1000G phase 1 and a local reference panel. They concluded that higher false positive rate was
103 observed in imputation from global reference panels compared to imputation performed using a
104 local panel. Other studies (Huang et al., 2015; Das et al., 2016) found imputation accuracy
105 increases with higher number of haplotypes, specifically for variants with $MAF \leq 0.5\%$. For
106 Hispanic populations, Nelson and colleagues (Nelson et al., 2016) compared imputation
107 performances with 1000G phase 1 (N=1,092) vs. 1000G phase 3 (N=2,504), concluding that
108 phase 3 improved accuracy for variants with $MAF < 1\%$ by . Further, Nagy and colleagues (Nagy
109 et al., 2017) showed that HRC reference panel provides new insight for novel variants
110 particularly for rare variants in a family-based Scottish study cohort. Aforementioned studies
111 highlighted the need of a larger sized reference panel to improve imputation quality. Herzig and
112 colleagues (Herzig et al., 2018) assessed tools for haplotype phasing and their impact on
113 imputation in a population isolate of Campora in southern Italy, and showed that SHAPEIT2,
114 SHAPEIT3 and EAGLE2 were highly accurate in phasing; MINIMAC3, IMPUTE4 and
115 IMPUTE2 were found to be reliable for imputation. Roshyara and colleagues (Roshyara et al.,
116 2014) compared MaCH-Admix, IMPUTE2, MACH, MACH-Minimac in different ethnicities by
117 evaluating accuracy of correctly imputed SNPs; MaCH-Minimac outperformed SHAPEIT-
118 IMPUTE2 in subsamples of different ethnic groups. These studies demonstrated how employed
119 imputation algorithm determines quality of inferred SNPs.

120

121 However, no study to our knowledge has evaluated reference panels in tandem with different
122 imputation algorithms to assess imputation quality of inferred SNPs based on MAF in a three-
123 way admixed population. Based on these findings, we assessed imputation quality, focusing on
124 rare and ultra-rare variants, in a large dataset of Caribbean Hispanics (CH) leveraging available
125 GWAS and sequencing data available for our cohort.

126

127

128

129 **Methods**

130 We will refer SNPs with MAF between 1-5% as “uncommon,” 0.1-1% as “rare,” and $\leq 0.1\%$ as
131 “ultra-rare”. We considered SNPs with IMPUTE-Info metric ≥ 0.40 as “good-quality” and ≥ 0.80
132 as “high-quality”.

133
134 GWAS samples and genotyping. We selected randomly 1,000 Caribbean Hispanics as part of an
135 original genotyped cohort of 3,138 individuals: **genotyped data can be downloaded at dbGaP**
136 **Study Accession: phs000496.v1.p1.** 719 individuals were derived from Estudio Familiar
137 Investigar Genetica de Alzheimer (EFIGA), a study of familial LOAD; and 281 individuals from
138 the multiethnic longitudinal cohort, Washington Heights, Inwood, Columbia Aging Project
139 (WHICAP). The information on study design, recruitment and GWAS methods for the EFIGA
140 and WHICAP study was previously described in Tosto, G., et al (Tosto et al., 2015).

141
142 GWAS quality control (QC). Genotyped data underwent quality control using PLINK (v1.90b4.9
143 64-bit) (Purcell et al., 2007). Briefly, we excluded SNPs with missing rate $\geq 5\%$ followed by
144 exclusion of SNPs with MAF $\leq 1\%$. We then removed SNPs with P-value $< 1e-6$ for Hardy-
145 Weinberg Equilibrium. Samples with missing call rate $\geq 5\%$ were excluded from analysis.

146
147 Global Ancestry estimation and selection of “true Hispanics”. Prior to imputation, we estimated
148 global ancestry using the ADMIXTURE (v.1.3.0) software (Alexander et al., 2009; Zhou et al.,
149 2011). We conducted supervised admixture analyses using three reference populations: African
150 Yoruba (YRI) and non-Hispanic white of European Ancestry (CEU) from the HAPMAP project
151 as representative of African and European ancestral populations; and eight Surui, 21 Maya, 14
152 Karitiana, 14 Pima and seven Colombian individuals from the Human Genome Diversity Project

153 (HGDP) were used to represent native American ancestry (Li et al., 2008). We used ~80,000
154 autosomal SNPs that were: I) genotyped in all three datasets (Caribbean Hispanics, 1000G and
155 HGDP); II) common (i.e. MAF >5 %); and III) in linkage equilibrium. Supervised admixture
156 analyses with the three reference populations (YRI, CEU, and Native Americans) revealed that
157 European lineage accounted for most of the ancestral origins (59%), followed by African (33%)
158 and native American ancestry (8%). We then selected only individuals with at least 1% of all
159 three ancestral populations.

160
161 *Reference panels.* HRC reference panel contained over 39M SNPs from 27,165 individuals who
162 participated in 17 different studies (**Table 1**). The data were downloaded from the Wellcome
163 Trust Sanger Institute (WTSI).

164 1000G phase 3 reference panel contained over 81M SNPs from 2,504 individuals
165 (https://mathgen.stats.ox.ac.uk/impute/1000GP_Phase3.tgz). It includes 26 ethnic groups, with
166 most variants rare, approximately 64 million had MAF <0.5%; approximately 12 million had a
167 MAF between 0.5% and 5%; and approximately 8 million have MAF >5%. In order to perform
168 imputation with MaCH-Admix, 1000G Phase 3 pre-formatted data were downloaded from
169 [ftp://yunlianon:anon@rc-ns-](ftp://yunlianon:anon@rc-ns-ftp.its.unc.edu/ALL.phase3_v5.shapeit2_mvncall_integrated.noSingleton.tgz)
170 ftp.its.unc.edu/ALL.phase3_v5.shapeit2_mvncall_integrated.noSingleton.tgz that contained over
171 47M SNPs.

172 The subsequent analyses were restricted to autosomal chromosomes, only.

173
174 *Phasing and Imputation Procedures.* We compared SHAPEIT2 (Delaneau et al., 2013) and Eagle2
175 (Loh et al., 2016) by phasing and then imputing (see next section) a single chromosome

176 (Chromosome 21), using both reference panels. We refer to SHAPEIT2 as SHAPEIT when used
177 in tandem with IMPUTE2 for the remainder of paper.
178
179 Imputation was carried out using two bioinformatics tools: IMPUTE2 (Howie et al., 2009) and
180 MaCH-Admix (Liu et al., 2013). For both, imputation quality ranged from 0 to 1, with 0
181 indicating complete uncertainty in imputed genotypes, and 1 indicating no uncertainty in
182 imputed genotypes.
183 IMPUTE2 (version 2.3.2). IMPUTE2 uses an MCMC algorithm to integrate over the space of
184 possible phase reconstructions for genotypes data. We conducted imputation in non-overlapping
185 1MB chunk regions; chunk coordinates were specified using the “-int” option. Other options
186 were used with default parameters (**Supplementary section 1**). Briefly, we used a default
187 250KB buffer region to avoid quality deterioration on the ends of chunk region. “-Ne” value as
188 2000 suggested for robust imputation which scales linkage disequilibrium and recombination
189 error rate. MaCH-Admix. We used MaCH-Admix because it uses a method based on IBS
190 matching in a piecewise manner. The method breaks genomic region under investigation into
191 small pieces and finds reference haplotypes that best represent every small piece, for each target
192 individual separately. MaCH-Admix imputes in three steps: phasing, estimation of model
193 parameter that includes error rate and recombination rate and lastly, haplotype-based imputation.
194 MaCH-Admix (version Beta 2.0.185) was run on default parameters of 30 rounds, 100 states (--
195 autoFlip flag). Details can be found in supplementary file (**section 1**). We initially compared
196 performance between MaCH-Admix and IMPUTE2 using the 1000G reference panel for
197 Chromosome 21 only. We then proceeded to impute all remaining chromosomes with the tool
198 that performed better.

199

200 Imputation Performance Metrics. IMPUTE2 uses “Info” parameter to report imputation quality
201 that measures relative statistical information about SNP allele frequency from imputed data. It
202 reflects the information in imputed genotypes relative to the information if only the allele
203 frequency were known. “Info” metric is used to filter poorly imputed SNPs from IMPUTE2 and
204 is reported for all imputed SNPs. In addition, IMPUTE2 uses an internal metric known as R^2 ,
205 reported for genotyped SNPs only: it measures squared correlation between genotyped SNPs and
206 the same SNPs that have been first masked internally and then imputed. MaCH-Admix uses *Rsq*
207 to report imputation quality. The R^2 metric is also known as variance ratio, calculated as
208 proportion of empirically observed variance (based on the imputation) to the expected binomial
209 variance $p(1-p)$, where p is the minor allele frequency. A threshold of 0.30 is recommended to
210 filter out poorly imputed SNPs.

211 Despite quality measures from IMPUTE2 and MaCH-Admix being highly correlated (Marchini
212 and Howie, 2010), we calculated a *r2hat* score to generate a single common metric to assess
213 imputation quality across the software (Hancock et al., 2012) (v109,
214 http://www.unc.edu/~yunmli/tgz/r2_hat.v109.tgz).

215 We compared performance of MaCH-Admix and SHAPEIT-IMPUTE2 by: a) Reporting raw
216 SNP counts based on quality (MaCH-Admix “Rsq” and IMPUTE2 “Info”); b) Comparing *r2hat*
217 for overlapping imputed SNPs from both tools; c) Conducting a Wilcoxon Signed-Rank Test (R
218 v3.4.2) on *r2hat* value of overlapping SNPs.

219

220 We compared performance of Eagle2 and SHAPEIT2 phasing tools in tandem with IMPUTE2 as
221 imputation tools across reference panels by: a) Comparing their respective IMPUTE2 R^2 : b)

222 Conducting a Wilcoxon Signed-Rank Test on R^2 value; c) Reporting raw counts of imputed
223 SNPs based on IMPUTE2 “Info” metric and stratified by MAF bins (e.g. common, rare, ultra-
224 rare).

225 In all comparisons, the MAFs are estimated from imputed data according to the reference panel
226 employed. We retained monomorphic SNPs in our analyses for several reasons. A monomorphic
227 SNP in one study might not be monomorphic in other cohorts. This has profound affects, for
228 example, when performing meta-analysis across different studies. In addition, monomorphic
229 SNPs provide information about MAF across studies. Without the information it is difficult to
230 tell, for instance, if a SNP is monomorphic or failed quality control in that study.

231
232 *Agreement between Imputed and Sequence data.* To further test the quality of imputation -
233 without relying on software’s internal metrics (i.e. “Info” and R^2) - we calculated genotyped
234 concordance between imputed and WES data using the VCF-compare tool (v0.1.14-12-
235 gcdb80b8) (Danecek et al., 2011). First, we converted posterior probabilities obtained from
236 imputation into genotype data using the PLINK software (v1.90b4.9) by applying a threshold of
237 0.9 (**supplementary section 1**), such that SNPs that failed on this criterion were left uncalled.
238 For example, an imputed SNP with $P(G=0,1,2) = (0.01, 0.9, 0.09)$ would be called as a '1'
239 (heterozygous), whereas an imputed SNP with $P(G=0,1,2) = (0.2, 0.6, 0.2)$ would be left
240 uncalled. We restricted the comparison to overlapping SNPs between HRC, 1000G reference
241 panels and whole-exome sequencing (WES) data for Chromosome 14 only, on SNPs with 0%
242 missingness (plink --missing flag) in WES data. We also assessed variants’ agreement according
243 to different MAF bins for “high-quality” (“Info” ≥ 0.8) SNPs. The output resulted in number of
244 variant “mismatches”, i.e. the count of allele not matching between imputed and sequenced

245 variants per individual. To measure interrater reliability we computed Cohen’s kappa coefficient
246 (McHugh, 2012) for both the reference panels against WES data. Kappa coefficient ≤ 0
247 indicates no agreement, 0.01–0.20 as none to slight, 0.21–0.40 as fair, 0.41– 0.60 as moderate,
248 0.61–0.80 as substantial, and 0.81–1.00 as almost perfect agreement.

249

250 Effects of Ancestry on Imputation Quality. To assess how ancestry affected imputation quality,
251 we conducted a Poisson regression using R. We used percentage of global ancestry (European
252 (CEU), Native (NAT) and African (YRI) as predictors, and total number of mismatches as the
253 outcome; analyses were restricted to “high-quality” SNPs, only.

254

255 Imputation of G206A Mutation in PSEN1. To evaluate imputation performance of a specific rare
256 variant, we examined a founder mutation, p.Gly206Ala (G206A - rs63750082) in the *PSEN1*
257 gene (PSEN1-G206A) (Athán et al., 2001; Lee et al., 2015). The PSEN1-G206A mutation is a
258 rare variant observed primarily in Puerto Ricans with familial early onset Alzheimer’s disease
259 (EOAD), but it is rare in Puerto Ricans and other populations with late-onset Alzheimer’s
260 disease (LOAD) (Arnold et al., 2013). The mutation was present in the 1000G phase 3 reference
261 panel with an allele frequency of 0.001, but was absent in the HRC reference panel. To verify
262 whether individuals who were found to carry the PSEN1-G206A mutation based on 1000G-
263 imputation, they were genotyped using the KASP genotyping technology by LGC genomics
264 (<https://www.lgcgroup.com>), which uses allele-specific PCR for SNP calling. Agreement
265 between imputed and genotype data for the PSEN1-G206A mutation was then assessed. We also
266 tested the effect on imputation quality based on different IMPUTE2-parameters settings, more
267 specifically by modifying the chunk size (i.e. 1MB vs. 5 MB).

268 **Results**

269

270 *Comparison of Phasing Tools: Eagle2 vs. SHAPEIT2.* To select the optimal tool for phasing, we
271 compared SHAPEIT2 with Eagle2 using Chromosome 21 with 13,066 genotyped SNPs by
272 performing subsequent imputation with IMPUTE2 on phased outputs, and using both reference
273 panels. We found SHAPEIT2 better than Eagle2 when evaluated based on mean R^2 and “Info”
274 metric using either the reference panels. For instance, using the 1000G, we observed higher
275 mean R^2 for data phased with SHAPEIT2 as compared to Eagle2 (0.92 vs. 0.91; Wilcoxon p-
276 value < 0.001). Similarly, when HRC panel was employed, mean R^2 of 0.89 was observed for
277 SHAPEIT2 against 0.85 for Eagle2 (Wilcoxon signed-rank test p-value < 0.001).

278 SNP count comparison details can be found in **Supplementary Table 1-2**. Regardless of the
279 reference panel employed, we observed higher percentage of “high-quality” rare and ultra-rare
280 SNPs for SHAPEIT-IMPUTE2 than Eagle2-IMPUTE2. For instance, 1000G-imputation
281 retrieved 51.02% of “high-quality” rare SNPs using SHAPEIT-IMPUTE2 vs. 48.38% with
282 Eagle2-IMPUTE2. Detailed comparisons for different MAF bins and quality threshold can be
283 found in **Supplementary Section 2**. Nevertheless, we found Eagle2 faster than SHAPEIT2 when
284 computation times were compared; for instance, with HRC Eagle2 was ~6 times faster than
285 SHAPEIT2 (**Supplementary Table 3**). We therefore imputed the remaining chromosomes on
286 phased output from SHAPEIT2. Comparison of phasing tools by assessing switch error rate was
287 beyond the scope of this paper due to limited resources, for e.g., availability of phased reference
288 panel for an admixed population.

289

290 MaCH-Admix vs. IMPUTE2. We found that SHAPEIT-IMPUTE2 performed better than MaCH-
291 Admix. For Chromosome 21, we imputed 1,104,648 and 646,594 SNPs for SHAPEIT-
292 IMPUTE2 and MaCH-Admix, respectively; 549,091 SNPs were overlapping. For SHAPEIT-
293 IMPUTE2 we observed 446,591 bi-allelic SNPs with “Info” ≥ 0.40 , in contrast with 598,943
294 SNPs with $R_{sq} \geq 0.30$ from MaCH-Admix (**Supplementary Table 4**). SNP counts for different
295 MAF bins based on platform-specific quality index can be found in **Supplementary Table**
296 **5**. When the two outputs were compared in terms of r^2_{hat} , SHAPEIT-IMPUTE2 showed a
297 higher average r^2_{hat} of 0.62 against 0.36 from MaCH-Admix (Wilcoxon signed-rank test p-value
298 < 0.001). Also, MaCH-Admix was 109 times slower than IMPUTE2. (**Supplementary Table 6**),
299 thus, comparison between different panels using MaCH-Admix were excluded due to limited
300 resources. For the remaining of this manuscript, we focused on imputation employing SHAPEIT-
301 IMPUTE2, only.

302
303 Comparison between HRC and 1000G using SHAPEIT-IMPUTE2. Using SHAPEIT-IMPUTE2,
304 we imputed 81,240,392 and 38,532,090 SNPs across all autosomal chromosomes with 1000G
305 and HRC reference panels, respectively (**Table 2**).

306 Overall, we observed slightly higher mean R^2 with 1000G than with HRC panel (0.94 vs. 0.92;
307 Wilcoxon p-value < 0.001). Nevertheless, when the analyses were restricted to only “good-” and
308 “high-quality” SNPs, HRC consistently performed better: 60.82% of HRC-imputed SNPs were
309 “good-quality” and 48.87% were “high-quality” (Wilcoxon signed-rank test p-value < 0.001). On
310 the contrary, 40.32% of 1000G imputed SNPs were “good-quality” and 30.11% were “high-
311 quality”.

312 Further, we evaluated performance for uncommon, rare and ultra-rare SNPs. For “good-” and
313 “high-quality” SNPs, HRC outperformed 1000G. For example, HRC panel produced 62.85% of
314 “high-quality” rare SNPs, whereas 1000G had 53.83% (**Table 3**). When average imputation
315 “Info” quality was evaluated, HRC-imputation again performed better than with 1000G
316 (Wilcoxon p-value < 0.001) (**Figure 1**).

317
318 Next, we restricted our analyses to *overlapping* SNPs across the two reference panels only, based
319 on their chromosome and position mapping, reference and non-reference alleles. For “good-” and
320 “high-quality” SNPs, imputation in both panels performed similarly (**Table 2**). When restricted
321 to uncommon, rare and ultra-rare SNPs, we observed higher percentage of “good-” and “high-
322 quality” SNPs for HRC panel as compared to 1000G reference panel (**Table 3**). For example,
323 7.44% of HRC-imputed ultra-rare SNPs were “good-quality” vs. 4.95% with the 1000G. 1.69%
324 of HRC-imputed ultra-rare SNPs were “high-quality” vs. 0.75% with the 1000G. Further,
325 Wilcoxon test on “Info” value of “high-quality” ultra-rare SNPs (2,972) again showed better
326 performances when HRC was employed vs. 1000G (P-value < 0.001). Complete list of counts
327 and percentages across reference panels, MAF bins and quality score can be found in **Table 3**.
328

329 *The case of G206A and the effect of chromosomal chunk size on imputation quality.* SNP
330 rs63750082 is absent from HRC panel therefore no imputation was achieved. Using 1000G
331 reference panel, 12 individuals were imputed as G206A carriers. SNP rs63750082 was imputed
332 with an IMPUTE2 “Info” score of 0.48 using 1MB as chromosomal region parameter. When we
333 increased the chunk size to 5MB, IMPUTE-Info score drastically improved to 0.94 (**Figure 2**).
334 Those patients labeled as mutation-carriers according to imputation were then genotyped: all 12

335 were confirmed to be G206A carriers, therefore achieving a perfect imputation prediction (100%
336 agreement) for that specific SNP.

337

338 Genotype Concordance and Kappa Coefficient. Out of the 1,000 individuals included in our
339 study, 262 had whole exome sequencing (WES) data available (Raghavan et al., 2018). We had
340 14,157 overlapping SNPs in WES, HRC and 1000G reference panels with 0% missingness in
341 WES data on Chromosome 14; SNPs imputed with each reference panel were compared against
342 WES data separately. When concordance was evaluated, HRC panel performed slightly poorer,
343 despite showing higher number of “high-quality” variants as compared to 1000G (**Table 4**).

344 Using 1000G, we observed 3,542 rare and 35 ultra-rare “high-quality” SNPs; across 262
345 samples, we counted 1,245 $((1,245/(3,542*262))*100= 0.13\%)$ and 10 (0.10%) mismatches for
346 rare and ultra-rare respectively. Using HRC, we retrieved 3,759 rare and 93 ultra-rare “high-
347 quality” variants; we observed 2,439 (0.24%) and 32 (0.13%) mismatches for rare and ultra-rare
348 variants, respectively. Details about pipeline can be found in **supplementary section 3**.

349 Next, we computed Cohen’s kappa coefficient (K) for 14,157 imputed SNPs common in WES
350 and the two reference panels. For both HRC and 1000G-imputation, we observed Kappa (K) of
351 ~ 0.99 for both rare and ultra-rare “high-quality” variants(**Table 4**). Details about pipeline can be
352 found in **supplementary section 4**.

353

354 Effects of Ancestry on Imputation Quality. We evaluated the effect of individual ancestral
355 component separately on SNP mismatches for Chromosome 14 on 262 individuals. For both
356 reference panels we found that higher African ancestry (YRI) was associated with higher number
357 of mismatches (**Supplementary table 7**). For instance, with 1000G reference panel, for rare
358 variants (“Info” ≥ 0.80), we observed an estimate of 1.46 (P-value <0.001) for YRI component

359 (indicating that for each unit increase in YRI ancestry, it results in 1.46 additional mismatches).
360 Details on confidence intervals and robust standard errors can be found in **supplementary file**
361 **(Table 7 and Section 5)**. We did not observe significant effect of ancestry on “high-quality”
362 ultra-rare variants in both panels.
363

364 **Discussion**

365 This study examined imputation performances in a cohort Caribbean Hispanics, focusing on
366 uncommon, rare and ultra-rare variant, by comparing different phasing and imputation tools, as
367 well as evaluating the effects of different reference panels. Overall, uncommon and rare variants
368 can be well imputed in this population, characterized by a unique genetic background. Caribbean
369 Hispanics are admixed with 59% of their genetic component from European, 32% African, and
370 8% Native American ancestry (Tosto et al., 2015). Due to their genetic makeup and unique
371 linkage disequilibrium patterns, admixed populations offer unique opportunity in studying
372 complex diseases. First, disease prevalence varies across ethnic groups (Igartua et al., 2015) and
373 certain admixed populations show higher incidence rates and prevalence (e.g. Alzheimer's
374 disease, diabetes etc.) or lower ones (e.g. multiple sclerosis). Second, variants that are ethnic-
375 specific may explain a higher prevalence of the disease of interest in admixed groups.

376

377 In the present study, we examined multiple parameters of imputation using the Caribbean
378 Hispanics population. First, we found that imputation using SHAPEIT-IMPUTE2 phasing
379 generated better results than Eagle2-IMPUTE2, and SHAPEIT-IMPUTE2 is superior to MaCH-
380 Admix in terms of imputation performances and process time.

381 Using SHAPEIT-IMPUTE2, 1000G SNPs outnumbered HRC panel because of the higher
382 number of SNPs included in the reference panel itself. However, when we restricted our analyses
383 to overlapping “good-” and “high-quality” SNPs (i.e. those variants that most likely would be
384 included in association analyses), HRC-imputation outperformed 1000G with higher. The
385 superior performance of HRC over 1000G was confirmed also when we focused on uncommon,

386 rare and ultra-rare SNPs only. Our findings confirm data in literature, i.e. reference panels with
387 higher number haplotypes perform better in different scenarios.

388 Additional investigations are needed in order to apply our findings to other admixed and non-
389 admixed populations.

390
391 Overall, higher quality of imputation for rare and ultra-rare variants was also confirmed when we
392 tested results against sequencing data. Finally, higher YRI global ancestry was found to
393 significantly impair SNP imputation, suggesting that imputation quality decreases with increased
394 African ancestry.

395
396 Lastly, SHAPEIT-IMPUTE2 with 1000G reference panel was successful in identifying G206A
397 mutation carriers. We also noticed that imputation quality drastically improved when imputation
398 was conducted using large (5MB) chunk size as compared to small (1MB) chunks. This seems to
399 contradict previous observation: Zhang et al (Zhang et al., 2011) studied the effect of window
400 size on imputation in an African-Americans . They concluded that window size of 1MB could be
401 considered acceptable. Possible explanations for these different results might be the more
402 complex admixture of CH compare to AA (three-way vs. two-way admixed) and a more
403 complex LD pattern for the G206A region. Ultimately, we recommend to consider a wider
404 window size to achieve high-quality imputation in specific variants that fail under default
405 settings.

406
407 This work has limitations. First, we could carry out the comparison between the two reference
408 panels restricting the analyses to overlapping variants only, limiting our observation to a subset

409 of the variants included in the 1000G panel. This is a result of the HRC composition, which is
410 composed by several studies and ended up including only a consensus number of variants.
411 Second, we tested the agreement between imputed and sequenced variants in a smaller subset of
412 individuals that had both GWAS and WES data available.

413 **Acknowledgments**

414 We thank the EFIGA study participants and the EFIGA research and support staff for their
415 contributions to this study. This study was supported by funding from the National Institute on
416 Aging [R21AG054832 (GT); 5R37AG015473 and RF1AG015473 (RM); R56 AG051876 and
417 R01 AG058918 (JL)] and the BrightFocus Foundation (A2015633S (JL)).

418

419

420 **Conflict of Interest Statement**

421 The authors declare no conflict of interests.

422

423 **References.**

- 424 Alexander, D.H., Novembre, J., and Lange, K. (2009). Fast model-based estimation of ancestry
425 in unrelated individuals. *Genome Res* 19(9), 1655-1664. doi: 10.1101/gr.094052.109.
- 426 Arnold, S.E., Vega, I.E., Karlawish, J.H., Wolk, D.A., Nunez, J., Negron, M., et al. (2013).
427 Frequency and clinicopathological characteristics of presenilin 1 Gly206Ala mutation in
428 Puerto Rican Hispanics with dementia. *J Alzheimers Dis* 33(4), 1089-1095. doi:
429 10.3233/JAD-2012-121570.
- 430 Athan, E.S., Williamson, J., Ciappa, A., Santana, V., Romas, S.N., Lee, J.H., et al. (2001). A
431 founder mutation in presenilin 1 causing early-onset Alzheimer disease in unrelated
432 Caribbean Hispanic families. *JAMA* 286(18), 2257-2263.
- 433 Browning, B.L., and Browning, S.R. (2009). A unified approach to genotype imputation and
434 haplotype-phase inference for large data sets of trios and unrelated individuals. *Am J*
435 *Hum Genet* 84(2), 210-223. doi: 10.1016/j.ajhg.2009.01.005.
- 436 Danecek, P., Auton, A., Abecasis, G., Albers, C.A., Banks, E., DePristo, M.A., et al. (2011). The
437 variant call format and VCFtools. *Bioinformatics* 27(15), 2156-2158. doi:
438 10.1093/bioinformatics/btr330.
- 439 Das, S., Forer, L., Schonherr, S., Sidore, C., Locke, A.E., Kwong, A., et al. (2016). Next-
440 generation genotype imputation service and methods. *Nat Genet* 48(10), 1284-1287. doi:
441 10.1038/ng.3656.
- 442 Delaneau, O., Zagury, J.F., and Marchini, J. (2013). Improved whole-chromosome phasing for
443 disease and population genetic studies. *Nat Methods* 10(1), 5-6. doi: 10.1038/nmeth.2307.
- 444 Genomes Project, C., Auton, A., Brooks, L.D., Durbin, R.M., Garrison, E.P., Kang, H.M., et al.
445 (2015). A global reference for human genetic variation. *Nature* 526(7571), 68-74. doi:
446 10.1038/nature15393.
- 447 Gibson, G. (2012). Rare and common variants: twenty arguments. *Nat Rev Genet* 13(2), 135-
448 145. doi: 10.1038/nrg3118.
- 449 Ha, N.T., Freytag, S., and Bickeboeller, H. (2014). Coverage and efficiency in current SNP
450 chips. *Eur J Hum Genet* 22(9), 1124-1130. doi: 10.1038/ejhg.2013.304.
- 451 Hancock, D.B., Levy, J.L., Gaddis, N.C., Bierut, L.J., Saccone, N.L., Page, G.P., et al. (2012).
452 Assessment of genotype imputation performance using 1000 Genomes in African
453 American studies. *PLoS One* 7(11), e50610. doi: 10.1371/journal.pone.0050610.
- 454 Herzig, A.F., Nutile, T., Babron, M.C., Ciullo, M., Bellenguez, C., and Leutenegger, A.L.
455 (2018). Strategies for phasing and imputation in a population isolate. *Genet Epidemiol*
456 42(2), 201-213. doi: 10.1002/gepi.22109.
- 457 Howie, B., Fuchsberger, C., Stephens, M., Marchini, J., and Abecasis, G.R. (2012). Fast and
458 accurate genotype imputation in genome-wide association studies through pre-phasing.
459 *Nat Genet* 44(8), 955-959. doi: 10.1038/ng.2354.
- 460 Howie, B.N., Donnelly, P., and Marchini, J. (2009). A flexible and accurate genotype imputation
461 method for the next generation of genome-wide association studies. *PLoS Genet* 5(6),
462 e1000529. doi: 10.1371/journal.pgen.1000529.
- 463 Huang, J., Howie, B., McCarthy, S., Memari, Y., Walter, K., Min, J.L., et al. (2015). Improved
464 imputation of low-frequency and rare variants using the UK10K haplotype reference
465 panel. *Nat Commun* 6, 8111. doi: 10.1038/ncomms9111.

- 466 Igartua, C., Myers, R.A., Mathias, R.A., Pino-Yanes, M., Eng, C., Graves, P.E., et al. (2015).
467 Ethnic-specific associations of rare and low-frequency DNA sequence variants with
468 asthma. *Nat Commun* 6, 5965. doi: 10.1038/ncomms6965.
- 469 Lee, J.H., Cheng, R., Vardarajan, B., Lantigua, R., Reyes-Dumeyer, D., Ortmann, W., et al.
470 (2015). Genetic Modifiers of Age at Onset in Carriers of the G206A Mutation in PSEN1
471 With Familial Alzheimer Disease Among Caribbean Hispanics. *JAMA Neurol* 72(9),
472 1043-1051. doi: 10.1001/jamaneurol.2015.1424.
- 473 Li, J.Z., Absher, D.M., Tang, H., Southwick, A.M., Casto, A.M., Ramachandran, S., et al.
474 (2008). Worldwide human relationships inferred from genome-wide patterns of variation.
475 *Science* 319(5866), 1100-1104. doi: 10.1126/science.1153717.
- 476 Liu, E.Y., Li, M., Wang, W., and Li, Y. (2013). MaCH-admix: genotype imputation for admixed
477 populations. *Genet Epidemiol* 37(1), 25-37. doi: 10.1002/gepi.21690.
- 478 Liu, Q., Cirulli, E.T., Han, Y., Yao, S., Liu, S., and Zhu, Q. (2015). Systematic assessment of
479 imputation performance using the 1000 Genomes reference panels. *Brief Bioinform*
480 16(4), 549-562. doi: 10.1093/bib/bbu035.
- 481 Loh, P.R., Danecek, P., Palamara, P.F., Fuchsberger, C., Y, A.R., H, K.F., et al. (2016).
482 Reference-based phasing using the Haplotype Reference Consortium panel. *Nat Genet*
483 48(11), 1443-1448. doi: 10.1038/ng.3679.
- 484 Marchini, J., and Howie, B. (2010). Genotype imputation for genome-wide association studies.
485 *Nat Rev Genet* 11(7), 499-511. doi: 10.1038/nrg2796.
- 486 McHugh, M.L. (2012). Interrater reliability: the kappa statistic. *Biochem Med (Zagreb)* 22(3),
487 276-282.
- 488 Nagy, R., Boutin, T.S., Marten, J., Huffman, J.E., Kerr, S.M., Campbell, A., et al. (2017).
489 Exploration of haplotype research consortium imputation for genome-wide association
490 studies in 20,032 Generation Scotland participants. *Genome Med* 9(1), 23. doi:
491 10.1186/s13073-017-0414-4.
- 492 Nelson, S.C., Stilp, A.M., Papanicolaou, G.J., Taylor, K.D., Rotter, J.I., Thornton, T.A., et al.
493 (2016). Improved imputation accuracy in Hispanic/Latino populations with larger and
494 more diverse reference panels: applications in the Hispanic Community Health
495 Study/Study of Latinos (HCHS/SOL). *Hum Mol Genet* 25(15), 3245-3254. doi:
496 10.1093/hmg/ddw174.
- 497 Pei, Y.F., Zhang, L., Li, J., and Deng, H.W. (2010). Analyses and comparison of imputation-
498 based association methods. *PLoS One* 5(5), e10827. doi: 10.1371/journal.pone.0010827.
- 499 Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M.A., Bender, D., et al. (2007).
500 PLINK: a tool set for whole-genome association and population-based linkage analyses.
501 *Am J Hum Genet* 81(3), 559-575. doi: 10.1086/519795.
- 502 Raghavan, N.S., Brickman, A.M., Andrews, H., Manly, J.J., Schupf, N., Lantigua, R., et al.
503 (2018). Whole-exome sequencing in 20,197 persons for rare variants in Alzheimer's
504 disease. *Ann Clin Transl Neurol* 5(7), 832-842. doi: 10.1002/acn3.582.
- 505 Roshvara, N.R., Kirsten, H., Horn, K., Ahnert, P., and Scholz, M. (2014). Impact of pre-
506 imputation SNP-filtering on genotype imputation results. *BMC Genet* 15, 88. doi:
507 10.1186/s12863-014-0088-5.
- 508 Surakka, I., Sarin, A.-P., Ruotsalainen, S.E., Durbin, R., Salomaa, V., Daly, M., et al. (2016).
509 The rate of false polymorphisms introduced when imputing genotypes from global
510 imputation panels. *bioRxiv*. doi: 10.1101/080770.

- 511 Tosto, G., Fu, H., Vardarajan, B.N., Lee, J.H., Cheng, R., Reyes-Dumeyer, D., et al. (2015). F-
512 box/LRR-repeat protein 7 is genetically associated with Alzheimer's disease. *Ann Clin*
513 *Transl Neurol* 2(8), 810-820. doi: 10.1002/acn3.223.
- 514 Verma, S.S., de Andrade, M., Tromp, G., Kuivaniemi, H., Pugh, E., Namjou-Khales, B., et al.
515 (2014). Imputation and quality control steps for combining multiple genome-wide
516 datasets. *Front Genet* 5, 370. doi: 10.3389/fgene.2014.00370.
- 517 Zhang, B., Zhi, D., Zhang, K., Gao, G., Limdi, N.N., and Liu, N. (2011). Practical Consideration
518 of Genotype Imputation: Sample Size, Window Size, Reference Choice, and Untyped
519 Rate. *Stat Interface* 4(3), 339-352.
- 520 Zheng, H.F., Ladouceur, M., Greenwood, C.M., and Richards, J.B. (2012). Effect of genome-
521 wide genotyping and reference panels on rare variants imputation. *J Genet Genomics*
522 39(10), 545-550. doi: 10.1016/j.jgg.2012.07.002.
- 523 Zheng, H.F., Rong, J.J., Liu, M., Han, F., Zhang, X.W., Richards, J.B., et al. (2015).
524 Performance of genotype imputation for low frequency and rare variants from the 1000
525 genomes. *PLoS One* 10(1), e0116487. doi: 10.1371/journal.pone.0116487.
- 526 Zhou, H., Alexander, D., and Lange, K. (2011). A quasi-Newton acceleration for high-
527 dimensional optimization algorithms. *Stat Comput* 21(2), 261-273. doi: 10.1007/s11222-
528 009-9166-3.
- 529

530

531 **Table 1:** SNP counts in HRC and 1000G reference panel.

Reference Panel	Individuals	Autosomal variants	Bi-allelic SNPs	Multi-allelic SNPs
1000G Phase 3	2,504	81,706,022	77,818,332	3,887,690
HRC	27,165*	39,131,600	39,131,600	NA

532 *For Chromosome 1, the number of individuals were 22,691

533

534 **Table 2:** Type of imputed SNPs across reference panels.

Reference Panel	Multi-allelic SNPs			Bi-allelic SNPs			Total SNPs		
	Total SNPs	Info* ≥ 0.40 (%)	Info ≥ 0.80 (%)	Total SNPs	Info ≥ 0.40 (%)	Info ≥ 0.80 (%)	Total SNPs	Info ≥ 0.40 (%)	Info ≥ 0.80 (%)
All SNPs									
1000G	3,319,815	2,586,342 (77.90)	2,061,295 (62.09)	77,920,577	31,423,926 (40.32)	23,468,086 (30.11)	81,240,392	31,423,926 (41.86)	25,529,381 (31.42)
HRC	NA	NA	NA	38,532,090	23,436,980 (60.82)	18,833,790 (48.87)	38,532,090	23,436,980 (60.82)	18,833,790 (48.79)
SNPs overlapping HRC & 1000G									
1000G	NA	NA	NA	30,090,251	22,631,112 (75.21)	18,408,585 (61.17)	30,090,251	22,631,112 (75.21)	18,408,585 (61.17)
HRC	NA	NA	NA	30,090,251	22,438,268 (74.56)	18,395,036 (61.13)	30,090,251	22,438,268 (74.56)	18,395,036 (61.13)

535

536

537

538

539 **Table 3:** SNP Counts for all Bi-allelic uncommon, rare and ultra-rare SNPs.

540

MAF	1000G			HRC		
	Info ≥ 0	Info ≥ 0.40 (%)	Info ≥ 0.80 (%)	Info ≥ 0	Info ≥ 0.40 (%)	Info ≥ 0.80 (%)
All SNPs						
[1% - 5%]	6,025,281	5,989,223 (98.90)	5,441,982 (90.31)	5,434,996	5,421,257 (99.84)	5,061,904 (93.13)
[0.1% - 1%)	20,249,058	16,881,286 (83.36)	10,901,789 (53.83)	11,780,671	10,931,924 (92.79)	7,404,808 (62.85)
[0 - 0.1%)	44,562,205	1,490,434 (3.34)	242,717 (0.544)	15,055,433	828,256 (5.50)	174,673 (1.16)
SNPs overlapping HRC & 1000G						
[1% - 5%]	5,624,956	5,604,308 (99.63)	5,148,285 (91.52)	5,396,207	5,385,364 (99.79)	5,037,187 (93.34)
[0.1% - 1%)	11,875,603	10,442,603 (87.93)	7,027,312 (59.17)	10,945,899	10,268,136 (93.80)	7,060,908 (64.50)
[0 - 0.1%)	6,314,479	312,967 (4.95)	47,614 (0.75)	7,519,807	560,043 (7.44)	127,423 (1.69)

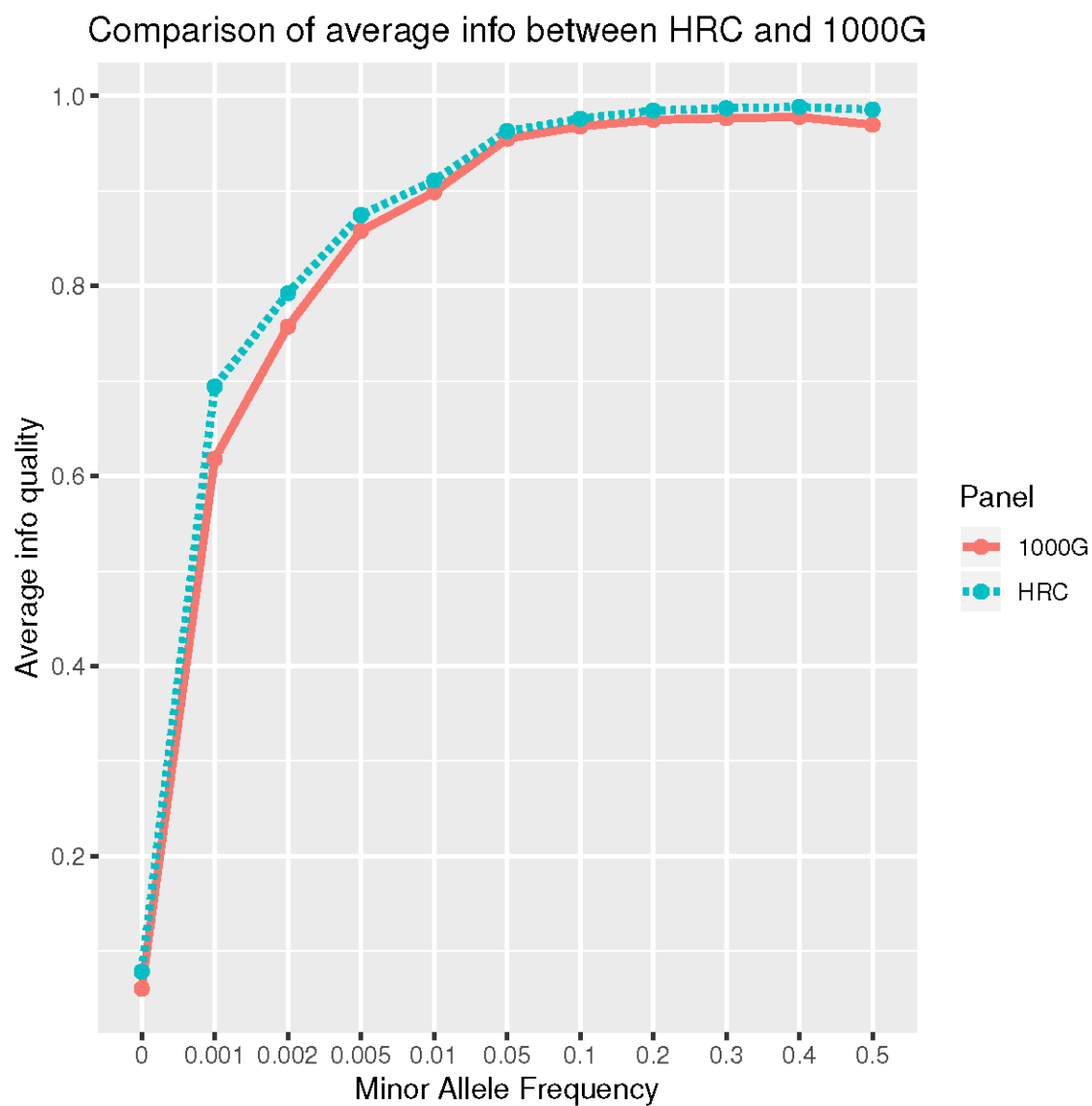
541

542 **Table 4:** Comparison for mismatch counts and Kappa (*K*) for HRC and 1000G using WES data
 543 on Chromosome 14.

MAF	1000G				HRC			
	Info ≥ 0.80				Info ≥ 0.80			
	SNP	Total SNPs in all persons*	Mismatch	Kappa (<i>K</i>)	SNP	Total SNPs in all persons*	Mismatch	Kappa (<i>K</i>)
[1% - 5%]	2,354	610,550	7,397 (1.22%)	0.99	2,264	587,961	8,963 (1.52%)	0.99
[0.1% - 1%]	3,542	926,109	1,245 (0.13%)	0.99	3,759	982,734	2,439 (0.24%)	0.99
[0 - 0.1%]	35	9,163	10 (0.10%)	0.99	93	24,348	32 (0.13%)	0.99

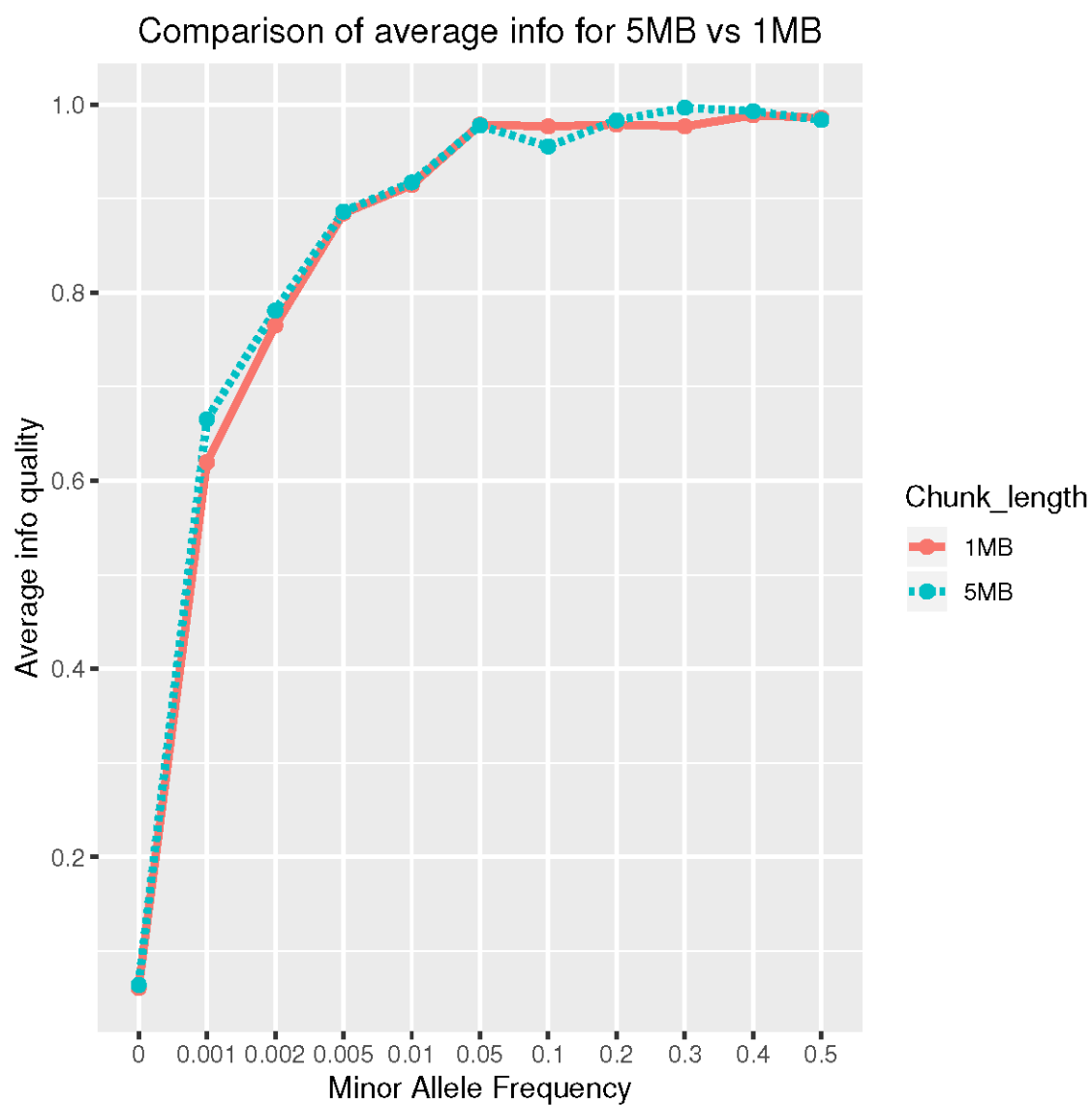
544 *Less value than 262*SNP because imputed with poor posterior probability failed to be
 545 converted from .gen to plink format.

546 **Figure 1:** Comparison of average Info quality between HRC and 1000G reference panel for all autosomal
547 chromosomes



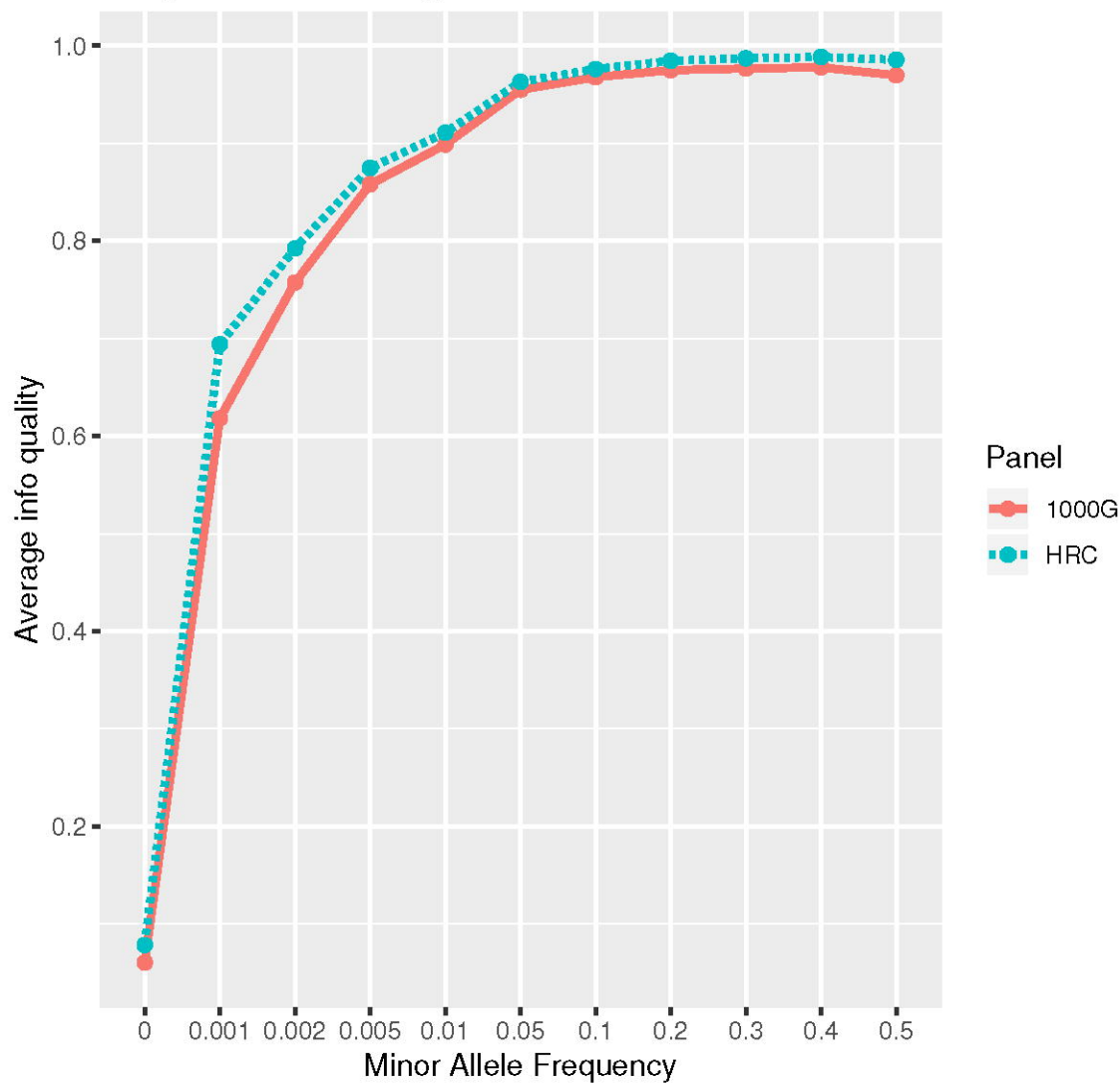
548

549 **Figure 2:** Comparison of average Info on CHR14: 70-75MB (5MB) vs 73-74MB (1MB) region



550

Comparison of average info between HRC and 1000G



Comparison of average info for 5MB vs 1MB

