

High-resolution population-specific recombination rates and their effect on phasing and genotype imputation

Running Title: Population-specific recombination maps in phasing & imputation

Shabbeer Hassan, *Institute for Molecular Medicine Finland, FIMM, HiLIFE, University of Helsinki, Helsinki, Finland*

Ida Surakka, *Institute for Molecular Medicine Finland, FIMM, HiLIFE, University of Helsinki, Helsinki, Finland*

Marja-Riitta Taskinen, *Clinical and molecular metabolism, Research program unit, University of Helsinki, Helsinki, Finland*

Veikko Salomaa, *Finnish Institute for Health and Welfare, Helsinki, Finland*

Aarno Palotie, *Institute for Molecular Medicine Finland, FIMM, HiLIFE, University of Helsinki, Helsinki, Finland, Psychiatric & Neurodevelopmental Genetics Unit, Department of Psychiatry, Analytic and Translational Genetics Unit, Department of Medicine, Department of Neurology, Massachusetts General Hospital, Boston, MA, USA*

Maija Wessman, *Institute for Molecular Medicine Finland, FIMM, HiLIFE, University of Helsinki, Helsinki, Finland*

Taru Tukiainen, *Institute for Molecular Medicine Finland, FIMM, HiLIFE, University of Helsinki, Helsinki, Finland*

Matti Pirinen, *Institute for Molecular Medicine Finland, FIMM, HiLIFE, University of Helsinki, Helsinki, Finland, Public Health, Clinicum, University of Helsinki, Helsinki, Finland, Department of Mathematics and Statistics, University of Helsinki, Helsinki, Finland*

Priit Palta, *Institute for Molecular Medicine Finland, FIMM, HiLIFE, University of Helsinki, Helsinki, Finland* , *Estonian Genome Center, Institute of Genomics, University of Tartu, Tartu, Estonia*

Samuli Ripatti, *Institute for Molecular Medicine Finland, FIMM, HiLIFE, University of Helsinki, Helsinki, Finland, Broad Institute of the Massachusetts Institute of Technology and Harvard University, Cambridge, MA, USA, Public Health, Clinicum, University of Helsinki, Helsinki, Finland*

CORRESPONDING AUTHOR

Samuli Ripatti, **Email:** samuli.ripatti@helsinki.fi

1 **Abstract:**

2 Recombination is an essential part of meiosis as it facilitates novel combinations of
3 homologous chromosomes, following their successive segregation in offspring. Founder
4 population size, demographic changes (eg. population bottlenecks or rapid expansion)
5 can lead to variation in recombination rates across different populations. Previous
6 research has shown that using population-specific reference panels has a significant
7 effect on downstream population genomic analysis like haplotype phasing, genotype
8 imputation and association, especially in the context of population isolates. Here, we
9 developed a high-resolution recombination rate mapping at 10kb and 50kb scale using
10 high-coverage (20-30x) whole-genome sequencing (WGS) data of 55 family trios from
11 Finland and compared it to recombination rates of non-Finnish Europeans (NFE). We
12 then tested the downstream effects of the population-specific recombination rates in
13 statistical phasing and genotype imputation in Finns as compared to the same analyses
14 performed by using the NFE-based recombination rates . Finnish recombination rates
15 have a moderately high correlation (Spearman's $\rho = 0.67-0.79$) with non-Finnish
16 Europeans, although on average (across all autosomal chromosomes), Finnish rates
17 (2.268 ± 0.4209 cM/Mb) are 12-14% lower than NFE (2.641 ± 0.5032 cM/Mb).
18 Population-specific effective population sizes were found to have no significant effect
19 in haplotype phasing accuracy (switch error rates, SER ~ 2%) and average imputation
20 concordance rates (with reference panels in phasing: rates were 97-98% for common,
21 92-96% for low frequency and 78-90% for rare variants) irrespective of the
22 recombination map used. Similarly, we found no effect of population-specific (Finnish)
23 recombination maps in phasing with comparable switch error rates (SER) across
24 autosomes when compared to HapMap based maps. Our results suggest that

25 downstream population genomic analyses like haplotype phasing and genotype
26 imputation mostly depend on population-specific contexts like appropriate reference
27 panels and their sample size, but not on population-specific recombination maps or
28 effective population sizes. Currently, available HapMap recombination maps seem
29 robust for population-specific phasing and imputation pipelines, even in the context of
30 relatively isolated populations like Finland.

31 Keywords: recombination, phasing, imputation, Finland, population genomics

32

33 **1. Introduction:**

34 Recombination is not uniform across the human genome with large areas having lower
35 recombination rates, so-called ‘coldspots’, which are then interspersed by shorter
36 regions marked by a high recombinational activity called ‘hotspots’ [1]. With long
37 chunks of human genome existing in high linkage disequilibrium, LD [2], and organised
38 in the form of ‘haplotype blocks’, the ‘coldspots’ tend to coincide with such regions of
39 high LD [3].

40 Direct estimation methods of recombination are quite time-consuming, and evidence
41 has suggested that they do not easily scale up to genome-wide, fine-scale
42 recombinational variation estimation [4]. A less time-consuming but computationally
43 intensive alternative is to use the LD patterns surrounding the SNPs [5]. Such methods
44 have been used in the past decade or so, to create fine-scale recombination maps [6].

45 Besides the International HapMap project that focused on capturing common variants
46 and haplotypes in diverse populations, international WGS-based collaborations like
47 1000 Genomes Project, provided genetic variation data for 20 worldwide populations

48 [7]. This led to further refinement of the recombination maps coupled with
49 methodological advances of using coalescent methods for recombination rate [8, 9].

50 With the rise of international collaborative projects, it was realised that founder
51 populations can often have very unique LD patterns [10], subsequently also displaying
52 unique increased genetics-driven health risks [11], suggesting that population-specific
53 reference datasets should be used to leverage the LD patterns to better detect disease
54 variants in downstream genetic analysis [12]. Genomic analysis methods like
55 haplotype phasing and imputing genotypes require recombination maps and other
56 population genetic parameters as input to obtain optimal results [13, 14, 15, 16]

57 In this study, we set to test this by 1) estimating recombination rates along the genome
58 in Finnish population using ~55 families of whole-genome sequenced (20-30x) Finns,
59 2) comparing these rates to some other European populations, and 3) comparing the
60 effect of using Finnish recombination rate estimates and cosmopolitan estimates in
61 phasing and imputation errors in Finnish samples.

62 **2. Materials & Methods:**

63 **2.1 Datasets used:**

64 *Finnish Migraine Families Collection*

65 Whole-genome sequenced trios (n = 55) consisting of the parent-offspring combination
66 were drawn from a large Finnish migraine families collection consisting of 1,589
67 families totalling 8,319 individuals [17]. The trios were used for the recombination map
68 construction using LDHAT version 2. The families were collected over 25 years from
69 various headache clinics in Finland (Helsinki, Turku, Jyväskylä, Tampere, Kemi, and
70 Kuopio) and via advertisements in the national migraine patient organisation web page
71 (<https://migreeni.org/>). The families consist of different pedigree sizes from small to

72 large (1-5+ individuals). Of the 8319 individuals, 5317 have a confirmed migraine
73 diagnosis based on the third edition of the established International Classification for
74 Headache Disorders (ICHD-3) criteria [18].

75 ***EUFAM cohort***

76 To check the phasing accuracy of our Finnish recombination map, we used an
77 independently sourced 49 trios from the European Multicenter Study on Familial
78 Dyslipidemias in Patients with Premature Coronary Heart Disease (EUfam). Finnish
79 familial combined hyperlipidemia (FCH) families were identified from patients initially
80 admitted to hospitals with premature cardiovascular heart disease (CHD) diagnosis who
81 also had elevated levels of total cholesterol (TC), triglycerides (TG) or both in the \geq
82 90th Finnish population percentile. Those families who had at least one additional first-
83 degree relative also affected with hyperlipidemia were also included in the study apart
84 from individuals with elevated levels of TG. [19, 20, 21].

85 ***FINRISK cohort***

86 The imputation accuracy of the Finnish and previously published HapMap based
87 recombination maps [8, 9] was subsequently tested on an independent FINRISK
88 CoreExome chip dataset consisting of 10,481 individuals derived from the national-
89 level FINRISK cohort. Primarily, it comprises of respondents of representative, cross-
90 sectional population surveys that are conducted once every 5 years since 1972 to get a
91 national assessment of various risk factors of chronic diseases and other health
92 behaviours among the working-age population drawn from 3 to 4 major cities in
93 Finland [22].

94 ***FINNISH reference panel cohort***

95 The whole-genome sequenced samples used were obtained from PCR-free methods and
96 PCR-amplified methods, which was followed by sequencing on a Illumina HiSeq X
97 platform with a mean depth of $\sim 30\times$. The obtained reads were then aligned to the
98 GRCh37 (hg19) human reference genome assembly using BWA-MEM. Best practice
99 guidelines from Genome Analysis Toolkit (GATK) were used to process the BAM files
100 and variant calling. Several criteria were used in this stage for sample exclusion:
101 relatedness (identity-by-descent (IBD) > 0.1), sex mismatches, among several others.
102 Furthermore, samples were filtered based on other criteria such as: non-reference
103 variants, singletons, heterozygous/homozygous variants ratio, insertion/deletion ratio
104 for novel indels, insertion/deletion ratio for indels observed in dbSNP, and
105 transition/transversion ratio.

106 After this stage, some exclusion criteria were applied to set some variants as missing:
107 $GQ < 20$, phred-scaled genotype likelihood of reference allele < 20 for heterozygous
108 and homozygous variant calls, and allele balance < 0.2 or > 0.8 for heterozygous calls. A
109 truth sensitivity percentage threshold of 99.8% for SNVs and of 99.9% for indels was
110 used based on the GATK Variant Quality Score Recalibration (VQSR) to filter variants
111 with, quality by depth (QoD) < 2 for SNVs and < 3 for indels, call rate $< 90\%$, and
112 Hardy-Weinberg equilibrium (HWE) p-value $< 1 \times 10^{-9}$. Some other variants like
113 monomorphic, multi-allelic and low-complexity regions [23] were further excluded.

114 The final reference dataset used in this study for imputation consisted of high coverage
115 (20-30x) whole-genome sequence-based reference panel of 2690 individuals from the
116 SISu project (Sequencing Initiative Suomi, <http://www.sisuproject.fi/>, [24]).

117 **2.2 Recombination map construction:**

118 Coalescent-based fine-scale recombination map construction [8] is greatly eased by
119 using trios which provide more accurate haplotype phasing resolution [25]. Hence, we
120 used trio data (n=55, 110 independent parents) from the Finnish Migraine Families
121 Cohort described above. These were filtered primarily using VCFtools [26] and custom
122 R scripts. Firstly, sites were thinned with within 15bp of each other such that only one
123 site remained followed by a filtering step of removing variants with a minor allele
124 frequency of <5% [27]. The resultant data were then phased using family-aware
125 method of SHAPEIT [28] using the standard HapMap recombination map [8, 9],
126 which was then split into segments of ~10000 SNPs with a 1000 SNP overhang on each
127 side of the segments. LDhat version 2 was run for 10^7 iterations with a block penalty of
128 5, every 5000 iterations of them of which the first 10% observations were discarded [8,
129 29]. The CEU based maps, used here for comparison, were obtained similarly using
130 LDhat [29].

131 However, LDhat is computationally intensive, and calculations suggest that the 1000
132 Genomes OMNI data set [30] would be too much computationally intensive to
133 complete [31], hence limiting the maximum number of haplotypes which could be
134 used.

135 To overcome this and make the recombination map independent of the underlying
136 methodology, we used a machine learning method implemented in FastEPRR [31, 32].
137 It supports the use of larger sample sizes, than LDhat and the recombination estimates
138 for sample sizes > 50, yields smaller variance than LDhat based estimates [31]. The
139 method was then applied to each autosome with overlapping sliding windows (*i.e.*,
140 window size, 50 kb and step length, 25 kb) under default settings for diploid organisms.

141 As seen in [31] both methods produce similar estimates, with only variance of the
142 estimate of mean being different.

143 The output of LDHat and FastEPRR is in terms of population recombination rate (p)
144 and to convert them into per-generational rate (r) used in phasing/imputation algorithms
145 we used optimal effective population size values derived from our testing (as explained
146 in the Supplementary Text). The estimates from LDHat and FastEPRR were then
147 averaged, to obtain a new combined estimate with the lowest variance amongst all the
148 three [31].

149 **2.3 Phasing and imputation accuracy**

150 To test whether the usage of different recombination maps affects the efficiency of
151 haplotype phasing and imputation, we used the aforesaid Finnish genotype data to
152 evaluate: (i) switch error rates across all chromosomes and (ii) imputation concordance
153 rates for chromosome 20.

154 **2.3.1 Phasing Accuracy**

155 The gold standard method to estimate haplotype phasing accuracy is to count the
156 number of switches (or recombination events) needed between the computationally
157 phased dataset and the true haplotypes [33]. The number of such switches divided by
158 the number of all possible switches is called switch error rate (SER).

159 For testing the influence of recombination maps on phasing accuracy, we used three
160 different recombination maps: HapMap, fine-scale Finnish recombination map and a
161 constant background recombination rate (1cM/Mb), to phase the 55 offspring
162 haplotypes without using any reference dataset. To check whether reference panels used
163 during haplotype phasing made any impact on the switch error rates, we used the

164 Finnish SISU based reference (n=2690), to check whether the size of the reference
165 panel made any impact on the results in phasing the offspring's haplotypes (Figure 1).
166 The SER in the offspring's phased haplotypes were then calculated by determining the
167 true offspring haplotypes using data from the parents (98 individuals) with a custom
168 script [34].

169 **2.3.2 Imputation Accuracy**

170 Imputation concordance was used as the metric for calculating the imputation accuracy.
171 For this, we randomly masked FINRISK CoreExome chip data consisting of 10,480
172 individuals [22] from chromosome 20. To test the role of reference panel size in
173 influencing the imputation accuracy in conjunction with varying the population genetics
174 parameters, we imputed the masked dataset with BEAGLE (Browning *et al.*, 2016)
175 using the Finnish reference panel (n = 2690). The concordance was then calculated
176 between the imputed genotypes and the original masked variants. Masking was done by
177 randomly removing ~10% of variants from the chip dataset.

178 The influence of recombination maps on imputation accuracy was checked by
179 calculating the concordance values between imputed and original variants, using the
180 Finnish reference panel in various combinations of recombination maps (constant rate,
181 HapMap, Finnish map) during the imputation (Figure 1).

182 **3. Results:**

183 **3.1 Finnish recombination map and its comparison to the HapMap recombination** 184 **map:**

185 The primary aim of our study was to derive a high-resolution genetic recombination
186 map for Finland and use it for comparative tests in commonly used analyses like
187 haplotype phasing and imputation. To derive a population-specific Finnish

188 recombination map, we used the high-coverage WGS data and an average of different
189 estimation methods (LDHat and FastEPRR). We used the N_e value of 10,000 derived
190 from our extensive testing of different N_e values (See supplementary text) to get the
191 per-generation recombination rates. The average recombination rates of Finnish
192 population isolate depicted 12-14% lower values (autosomal-wide average 2.268 ± 0.4209
193 cM/Mb) for all chromosomes compared to CEU based maps (2.641 ± 0.5032 cM/Mb)
194 (Figure 2).

195 These differences in average recombination rates are reflected in the correlation values
196 across all chromosomes (Spearman's $\rho \sim 0.67 - 0.79$) between the developed Finnish
197 map and HapMap based one (Figure 2). We also present a direct comparison between
198 the two maps, of the recombination rates at 5Mb scales, which presents a similar visual
199 pattern of rates across the genome (Supplementary Figure 1).

200 **3.2 Effects of the population-specific recombinations map on haplotype phasing**

201 Variation in population-specific recombination maps (and effective population sizes)
202 can affect the downstream genomic analyses like haplotype phasing and imputation.

203 We tested the Finnish map, HapMap map and a constant recombination rate map
204 (1cM/Mb) to understand the effects of population-specific maps on downstream
205 genomic analyses. The phasing accuracy was tested under two different conditions:
206 using no additional reference panel and using an population-specific .SISu v2 reference
207 panel ($n = 2690$) in phasing. We observed that, on average, SER ranged between 1.8-
208 3.7% (Supplementary Figure 2) across the different chromosomes and recombination
209 maps. We found statistically significant differences within both no-reference panel and
210 the Finnish reference panel results (Kruskal Wallis, p -value = $5.3e-10$ and $4.7e-10$,
211 respectively; Figure 3). The constant recombination map (1cM/Mb) had significantly

212 higher SER values when compared to the Finnish map or the HapMap map (Figure 3)
213 both when no reference panels were used (p-value = $2.9e-11$ and $2.6e-09$, respectively)
214 and when the Finnish reference panel was used (p-value = $2.9e-11$ and $9.5e-13$,
215 respectively). The choice of recombination maps mattered more when no reference
216 panel was used (p-value = 0.0046), however when using the Finnish reference panel, the
217 difference in SER was statistically insignificant (p-value = 0.25).

218 **3.3 Effects of the population-specific recombinations map on genotype imputation**

219 Imputation accuracy was similarly tested using the reference panel under three different
220 recombination map settings. We observed that when the imputation target dataset was
221 phased and imputed using the Finnish reference panel (n=2690) irrespective of the
222 population-specific recombination maps, it had a high imputation accuracy (overall
223 concordance rate ~98%, Figure 4) across MAF bins ($>0.1\%$). Though some differences
224 in concordance rates are seen in for rare variants (MAF $<0.1\%$). The concordance rate
225 was lower when the test dataset was phased without reference panels (concordance rate
226 72~77%, Figure 5).

227 **4. Discussion:**

228 Population isolates like Finland, have had a divergent demographic history as compared
229 to the outbred European populations, with a less historic migration, more fluctuating
230 population sizes and higher incidences of bottleneck events and founder effects [35, 36]
231 This unique demographic history then affects different population genetic parameters,
232 like recombination rates [37]. It has been shown previously that using population-
233 specific genomic reference panels augmented the accuracy of imputation accuracy
234 leading to better mapping of diseases specific variants in GWAS [12]. Since
235 recombination rates (in the form of recombination maps), features in much of the

236 downstream genomic analyses' methods like imputation and haplotype phasing [15,
237 34], we wanted to study their effect on downstream analyses.

238 Firstly, we characterised the Finnish recombination map using high-coverage (~30x)
239 whole-genome sequencing (WGS) samples from large SISu v2 reference panel
240 (n=2690). Previously used recombination maps hail from the HapMap and
241 1000Genomes projects which used sparse genotypic datasets or low-depth sequencing
242 samples. This is a first attempt in creating a recombination map for Finland using
243 population-specific WGS samples. We used two different methods in estimating the
244 recombination rates, to achieve accurate estimates with lower variance [29,31]. In
245 addition, we estimated effective population sizes using identity-by-descent (IBD) based
246 methods [15] for both Finnish and CEU based datasets. The obtained recombination
247 map was then used to test their role and importance in two selected downstream
248 genomic analyses – haplotype phasing and imputation concordance. Since the
249 recombination rate determination requires effective population size estimates, we also
250 tested the role of varying effective population size on these two analyses (See
251 Supplementary Text). The extensive testing of N_e yielded the estimate of 10,000
252 originally derived theoretically [38] and most used commonly for humans fits quite
253 rightly for the recombination map.

254 The Finnish recombinational landscape when compared to the HapMap based map,
255 showed, on average, a high degree of correlation across scales (10, 50kb and 5Mb),
256 however, on average, Finnish recombination rates across chromosomes were found to
257 be lower. Such moderate to high correlations (Figure 2) and similar recombinational
258 landscape (Supplementary Figure 1) could be due to high sharing of recombinations in
259 individuals from closely-related populations. The degree of dissimilarity in the

260 population-level differences between Finnish and mainland Europeans in terms of
261 recombination rates could be due to population-specific demographic processes like
262 founder effects, bottleneck events and migration [39], or chromatin structure PRDM9
263 binding locations, for example [40]. And the broad similarity in terms of correlational
264 structure seen here, reflects a shared ancestral origin of Finns and other mainland
265 Europeans [41]. Other studies on population isolates like Iceland [9] have previously
266 found a high degree of correlation with CEU based maps, albeit with substantial
267 differences as seen here. Previous studies [42] have additionally explored the
268 relationship between recombination rate differences between populations and allele
269 frequency differences, with evidence suggesting that the differences between rates show
270 the selection impact in the past 100,000 years since the out-of-Africa movement of
271 humans.

272 As seen in previous studies, much of the downstream genomic analyses like getting
273 more refined GWAS hits or, accurate copy number variants (CNV) imputation, can be
274 highly improved with the addition/use of population-specific datasets [12]. To test this
275 in the context of population-specific recombination maps, we used them to test the
276 haplotype phasing and imputation accuracy and observed that despite large differences
277 in the effective population sizes between populations, it did not affect the tested metrics.
278 One possible explanation for the insignificant effect seen here is that the role of
279 parameters like effective population size and recombination maps is to scale over the
280 haplotypes for efficient coverage of the whole genome. However, when sufficiently
281 large, population-specific genomic reference panels are available with tens of thousands
282 of haplotypic combinations, such scaling over for specific populations, does not yield
283 in substantial improvements. As we showed here, reference panel size could play an

284 important role in the downstream genomic analyses and in most cases, the current
285 practice of using the standard HapMap recombination map can be reasonably used.
286 Another point of interest here is that the use of different N_e parameters during
287 phasing/imputation might be redundant as we observed no change in the accuracy of our
288 estimates on varying the N_e parameters. Similarly, when using population-specific
289 recombination maps, we did not find any tangible benefits in using them over the
290 current standard maps based on the HapMap data.

291 Our study suggests a couple of important points for future studies: (a) varying effective
292 population size for downstream genomic analyses, such as phasing and imputation,
293 might have a relatively small impact, and it might be better to use the default option of
294 the particular software; (b) when available, it is beneficial to use a population-specific
295 genomic reference panel as they increase the accuracy; (c) HapMap can be used for
296 current downstream genomic analyses like haplotype phasing or genotype imputation in
297 European-based populations. And, if need be, can be substituted for using population-
298 specific maps, as the accuracy rates are quite similar to the population-based maps.

299 Though the sample used here is from a disease cohort but is nevertheless representative
300 of Finland's population and hence provides a reasonable recombination rate estimates.

301 On the other hand, our reliance on disease cohorts could lead to minor variation in the
302 resultant recombination. Though as we share a similar out-of-Africa origin, much of our
303 history is shared and though biological differences in the recombinational landscape do
304 exist between different populations, much of the downstream genomic analyses
305 (haplotyping, imputation or, GWAS), might not be affected by recombination map or
306 values of effective population size.

307 **Funding**

308 This work was financially supported by the Academy of Finland (251217 and 255847 to
309 S.R.). S.R. was further supported by the Academy of Finland Center of Excellence for
310 Complex Disease Genetics, the Finnish Foundation for Cardiovascular Research,
311 Biocentrum Helsinki, and the Sigrid Jusélius Foundation. S.H. was supported by
312 FIMM-EMBL PhD program doctoral funding and I.S. by Academy of Finland
313 Postdoctoral Fellowship (298149). V.S. was supported by the Finnish Foundation for
314 Cardiovascular Research. T.T. was supported by Academy of Finland grant number
315 315589.

316 **Acknowledgements**

317 We would like to thank Sari Kivikko and Huei-Yi Shen for management assistance. The
318 FINRISK analyses were conducted using the THL biobank permission for project
319 BB2015_55.1. We thank all study participants for their generous participation in the
320 FINRISK study.

321 *Conflict of Interest:* VS has received honoraria from Novo Nordisk and Sanofi for
322 consulting and has ongoing research collaboration with Bayer ltd (all unrelated to the
323 present study).

324 **References**

- 325 1. Baudat F, Buard J, Grey C, Fledel-Alon A, Ober C, Przeworski M et al. PRDM9
326 is a major determinant of meiotic recombination hotspots in humans and mice.
327 Science 2009; 327: 836–840
- 328 2. Daly MJ, Rioux JD, Schaffner SF, Hudson TJ, Lander ES. High-resolution
329 haplotype structure in the human genome. Nature Genetics 2001; 29: 229–232

- 330 3. Hudson RR, Kaplan NL. Statistical properties of the number of recombination
331 events in the history of a sample of DNA sequences. *Genetics* 1985; 111: 147-
332 164
- 333 4. Chan AH, Jenkins PA, Song YS. Genome-wide fine-scale recombination rate
334 variation in *Drosophila melanogaster*. *PLoS Genet* 2012; 8: e1003090
- 335 5. McVean GA, Myers SR, Hunt S, Deloukas P, Bentley DR, Donnelly P. The
336 fine-scale structure of recombination rate variation in the human genome.
337 *Science* 2004; 304: 581-584
- 338 6. Myers S, Bottolo L, Freeman C, McVean G, Donnelly P. A fine-scale map of
339 recombination rates and hotspots across the human genome. *Science* 2005; 310:
340 321-324.
- 341 7. Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO et al. A
342 global reference for human genetic variation. *Nature* 2015; 526: 68-74
- 343 8. Auton A, McVean G. Recombination rate estimation in the presence of hotspots.
344 *Genome Res* 2007; 17: 1219-1227.
- 345 9. Kong A, Thorleifsson G, Gudbjartsson DF, Masson G, Sigurdsson A,
346 Jonasdottir A et al. Fine-scale recombination rate differences between sexes,
347 populations and individuals. *Nature* 2010; 467: 1099-1103.
- 348 10. Service S, DeYoung J, Karayiorgou M, Roos JL, Pretorius H, Bedoya G et al.
349 Magnitude and distribution of linkage disequilibrium in population isolates and
350 implications for genome-wide association studies. *Nat Genet* 2006; 38: 556-560.
- 351 11. Peltonen L, Jalanko A, Varilo T. Molecular genetics of the Finnish disease
352 heritage. *Hum Mol Genet* 1999; 8: 1913-1923.

- 353 12. Surakka I, Kristiansson K, Anttila V, Inouye M, Barnes C, Moutsianas L et al.
354 Founder population-specific HapMap panel increases power in GWA studies
355 through improved imputation accuracy and CNV tagging. *Genome Res* 2010;
356 20: 1344-1351.
- 357 13. Tewhey R, Bansal V, Torkamani A, Topol EJ, Schork NJ. The importance of
358 phase information for human genomics. *Nat Rev Genet* 2011; 12: 215-223.
- 359 14. Browning SR, Browning BL. Haplotype phasing: existing methods and new
360 developments. *Nat Rev Genet* 2011; 12: 703-714.
- 361 15. Browning BL, Browning SR. Genotype Imputation with Millions of Reference
362 Samples. *Am J Hum Genet* 2016; 98: 116-126.
- 363 16. Delaneau O, Zagury JF, Marchini J. Improved whole-chromosome phasing for
364 disease and population genetic studies. *Nat Methods* 2013; 10: 5-6.
- 365 17. Gormley P, Kurki MI, Hiekkala ME, Veerapen K, Häppölä P, Mitchell AA et al.
366 Common Variant Burden Contributes to the Familial Aggregation of Migraine
367 in 1,589 Families. *Neuron* 2018; 98: 743-753.e4.
- 368 18. The International Classification of Headache Disorders, 3rd edition (beta
369 version). *Cephalalgia* 2013; 33: 629-808.
- 370 19. Borodulin K, Vartiainen E, Peltonen M, Jousilahti P, Juolevi A, Laatikainen T et
371 al. Forty-year trends in cardiovascular risk factors in Finland. *Eur J Public*
372 *Health* 2015; 25: 539-546.
- 373 20. Porkka KV, Nuotio I, Pajukanta P, Ehnholm C, Suurinkeroinen L, Syväne M et
374 al. Phenotype expression in familial combined hyperlipidemia. *Atherosclerosis*
375 1997; 133: 245-253.

- 376 21. Ripatti P, Rämö JT, Söderlund S, Surakka I, Matikainen N, Pirinen M et al. The
377 Contribution of GWAS Loci in Familial Dyslipidemias. *PLOS Genetics* 2016;
378 12: e1006078.
- 379 22. Vartiainen E, Laatikainen T, Peltonen M, Juolevi A, Mannisto S, Sundvall J et
380 al. Thirty-five-year trends in cardiovascular risk factors in Finland. *International*
381 *Journal of Epidemiology* 2009; 39: 504–518.
- 382 23. Li H. Toward better understanding of artifacts in variant calling from high-
383 coverage samples. *Bioinformatics* 2014; 30: 2843–2851.
- 384 24. Mart Kals, Tiit Nikopensius, Kristi Läll, Kalle Pärn, Timo Tõnis Sikka, Jaana
385 Suvisaari, Veikko Salomaa, Samuli Ripatti, Aarno Palotie, Andres Metspalu,
386 Tõnu Esko, Priit Palta, Reedik Mägi Advantages of genotype imputation with
387 ethnically matched reference panel for rare variant association analyses bioRxiv
388 579201; doi: <https://doi.org/10.1101/579201>
- 389 25. Roach JC, Glusman G, Hubley R, Montsaroff SZ, Holloway AK, Mauldin DE et
390 al. Chromosomal haplotypes by genetic phasing of human families. *Am J Hum*
391 *Genet* 2011; 89: 382-397.
- 392 26. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA et al. The
393 variant call format and VCFtools. *Bioinformatics* 2011; 27: 2156-2158.
- 394 27. Stevison LS, Woerner AE, Kidd JM, Kelley JL, Veeramah KR, McManus KF et
395 al. The Time Scale of Recombination Rate Evolution in Great Apes. *Mol Biol*
396 *Evol* 2016; 33: 928-945.
- 397 28. O'Connell J, Gurdasani D, Delaneau O, Pirastu N, Ulivi S, Cocca M et al. A
398 general approach for haplotype phasing across the full spectrum of relatedness.
399 *PLoS Genet* 2014; 10: e1004234.

- 400 29. Auton A, Fledel-Alon A, Pfeifer S, Venn O, Séguirel L, Street T et al. A fine-
401 scale chimpanzee genetic map from population sequencing. *Science* 2012; 336:
402 193-198.
- 403 30. Abecasis GR, Auton A, Brooks LD, DePristo MA, Durbin RM, Handsaker RE
404 et al. An integrated map of genetic variation from 1,092 human genomes. *Nature*
405 2012; 491: 56-65.
- 406 31. Gao F, Ming C, Hu W, Li H. New Software for the Fast Estimation of
407 Population Recombination Rates (FastEPRR) in the Genomic Era. *G3*
408 (Bethesda) 2016; 6: 1563-1571.
- 409 32. Lin K, Futschik A, Li H. A fast estimate for the population recombination rate
410 based on regression. *Genetics* 2013; 194: 473-484.
- 411 33. Bansal V. Integrating read-based and population-based phasing for dense and
412 accurate haplotyping of individual genomes. *Bioinformatics* 2019; 35: i242-
413 i248.
- 414 34. Loh PR, Danecek P, Palamara PF, Fuchsberger C, A Reshef Y, K Finucane H et
415 al. Reference-based phasing using the Haplotype Reference Consortium panel.
416 *Nat Genet* 2016; 48: 1443-1448.
- 417 35. Martin AR, Karczewski KJ, Kerminen S, Kurki MI, Sarin AP, Artomov M et al.
418 Haplotype Sharing Provides Insights into Fine-Scale Population History and
419 Disease in Finland. *Am J Hum Genet* 2018; 102: 760-775.
- 420 36. Kerminen S, Havulinna AS, Hellenthal G, Martin AR, Sarin AP, Perola M et al.
421 Fine-Scale Genetic Structure in Finland. *G3 (Bethesda)* 2017; 7: 3459-3468.
- 422 37. Wang J, Santiago E, Caballero A. Prediction and estimation of effective
423 population size. *Heredity (Edinb)* 2016; 117: 193-206.

- 424 38. Takahata N, Satta Y, Klein J. Divergence time and population size in the lineage
425 leading to modern humans. *Theor Popul Biol* 1995; 48: 198-221.
- 426 39. Novembre J, Johnson T, Bryc K, Kutalik Z, Boyko AR, Auton A et al. Genes
427 mirror geography within Europe. *Nature* 2008; 456: 98-101.
- 428 40. Ségurel L. The complex binding of PRDM9. *Genome Biol* 2013; 14: 112.
- 429 41. Mallick S, Li H, Lipson M, Mathieson I, Gymrek M, Racimo F et al. The
430 Simons Genome Diversity Project: 300 genomes from 142 diverse populations.
431 *Nature* 2016; 538: 201-206.
- 432 Keinan A, Reich D. Human population differentiation is strongly correlated with
433 local recombination rate. *PLoS Genet* 2010; 6: e1000886.

434

435

436

437

438

439

440

441 **Figure 1:** Flowchart overview of the analyses and comparisons performed.

442 **Figure 2:** Average (\pm standard deviation) recombination rates of Finnish v/s CEU per
443 autosome measured in cM/Mb and Correlation between Finnish and CEU
444 recombination rates across all chromosomes. The comparisons are made for similar
445 physical positions.

446 **Figure 3:** Statistical comparison of Switch Error Rates across all autosomes calculated
447 for all children in the trios using different recombination maps with respect to different

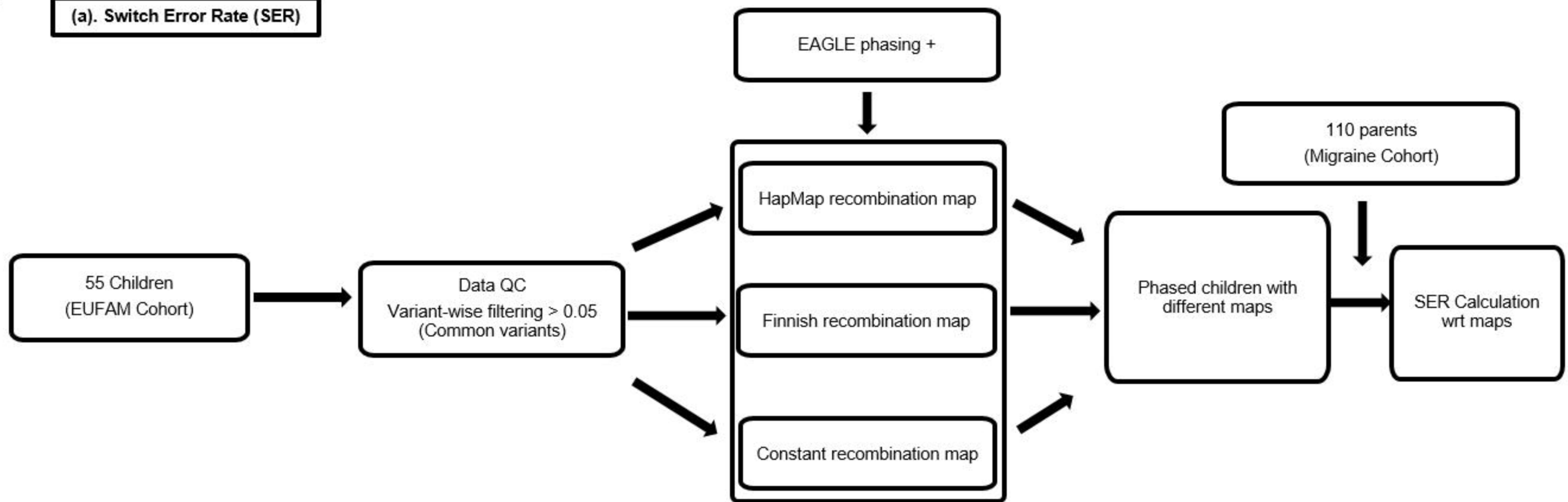
448 reference panel conditions (absent or present). The p-values are shown at the top of each
449 panel from Kruskal Wallis ANOVA testing between panel groups and ones between
450 boxplots for within-group comparisons.

451 **Figure 4:** Comparison of Imputation Concordance across different Minor Allele
452 Frequency (MAF) groups for a range of different recombination map combinations
453 phased with NO reference panels

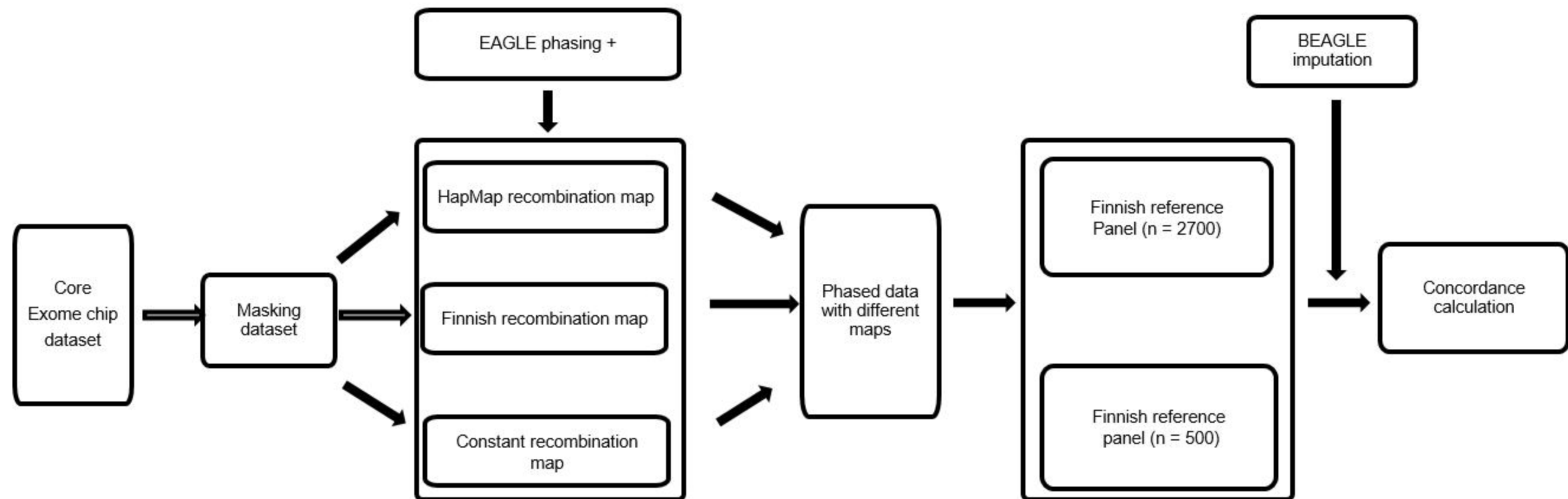
454 **Figure 5:** Comparison of Imputation Concordance across different Minor Allele
455 Frequency (MAF) groups for a range of different recombination map combinations
456 phased with reference panels.

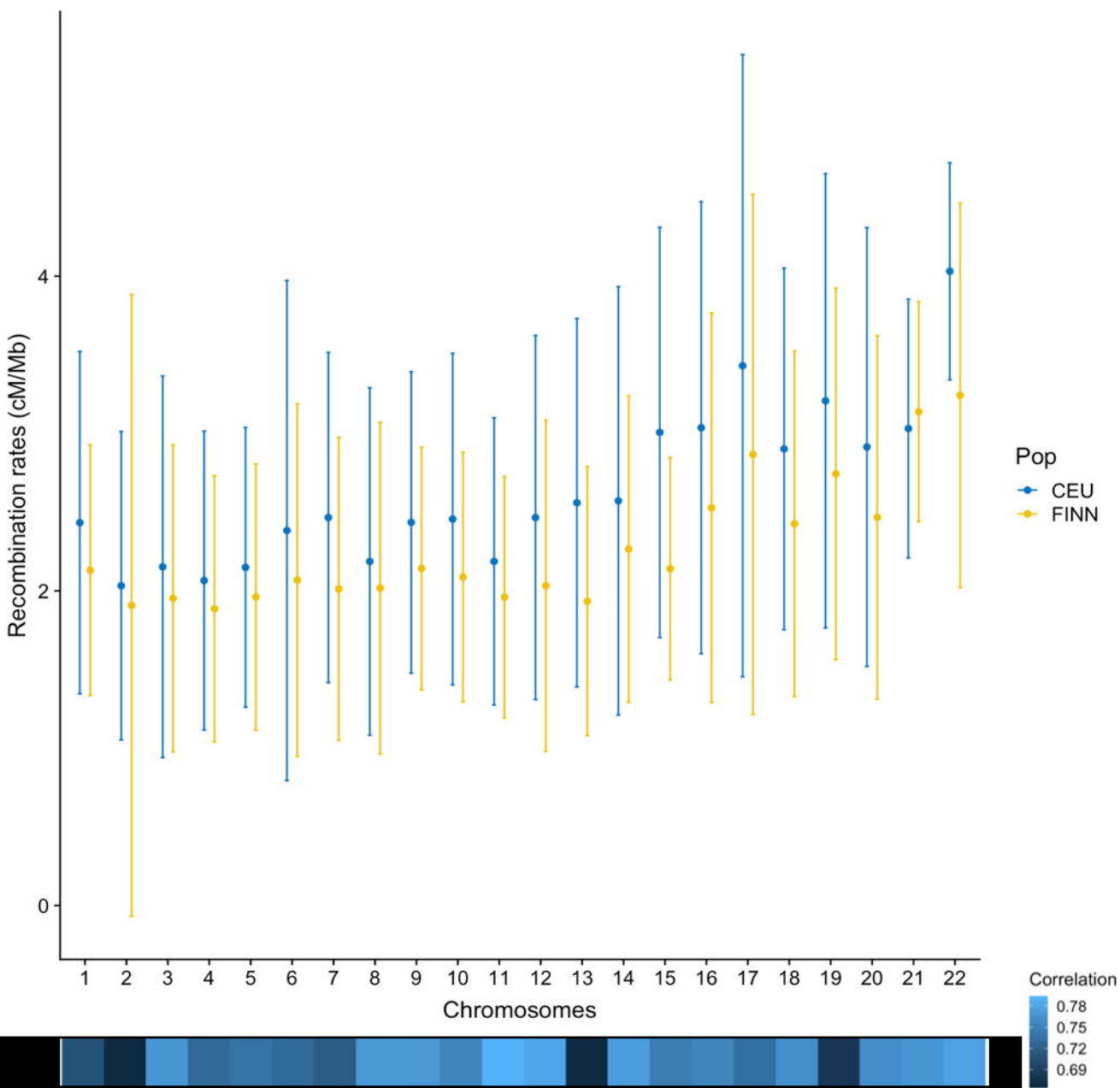
Recombination map comparison tests

(a). Switch Error Rate (SER)

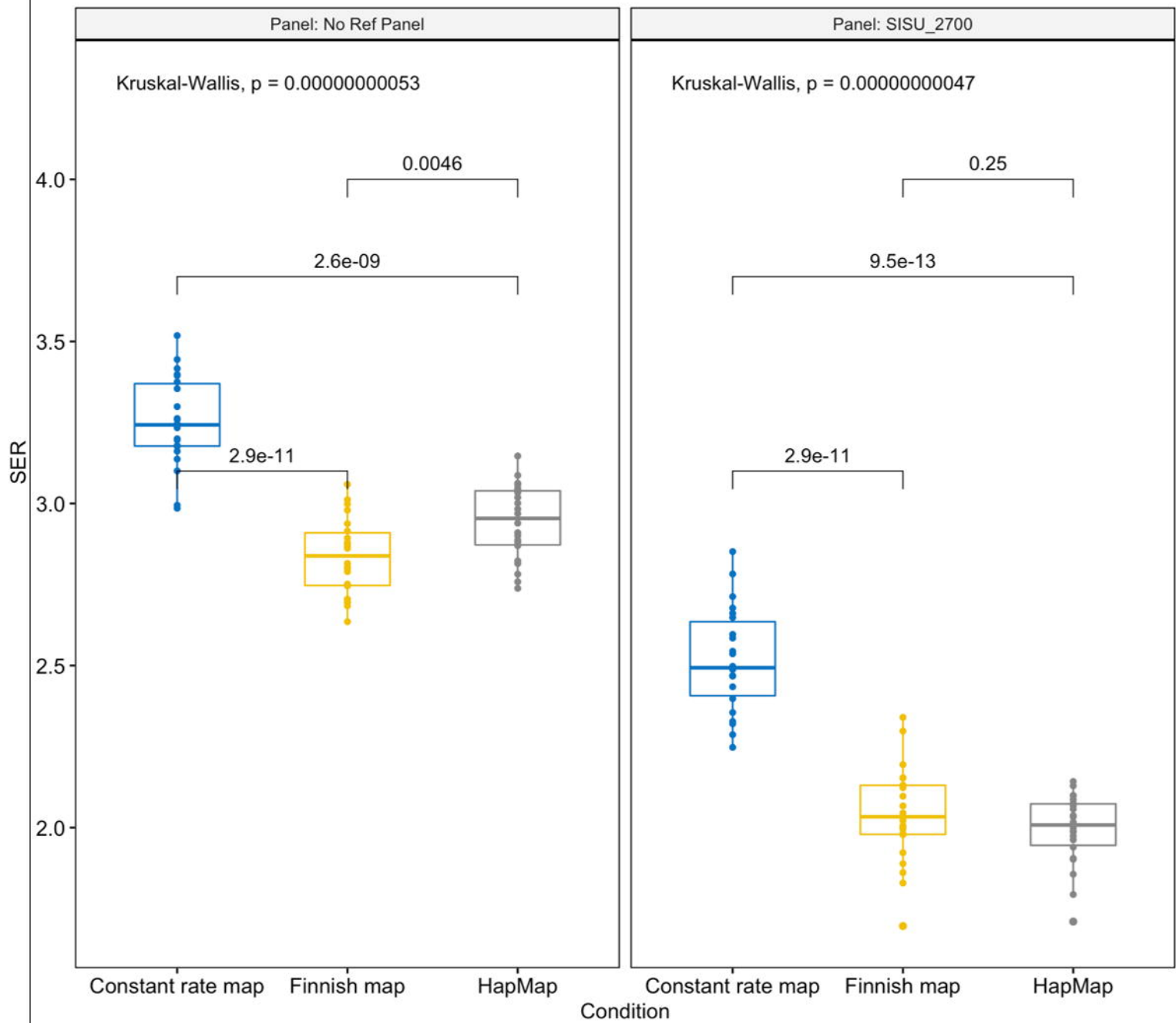


(b). Imputation concordance

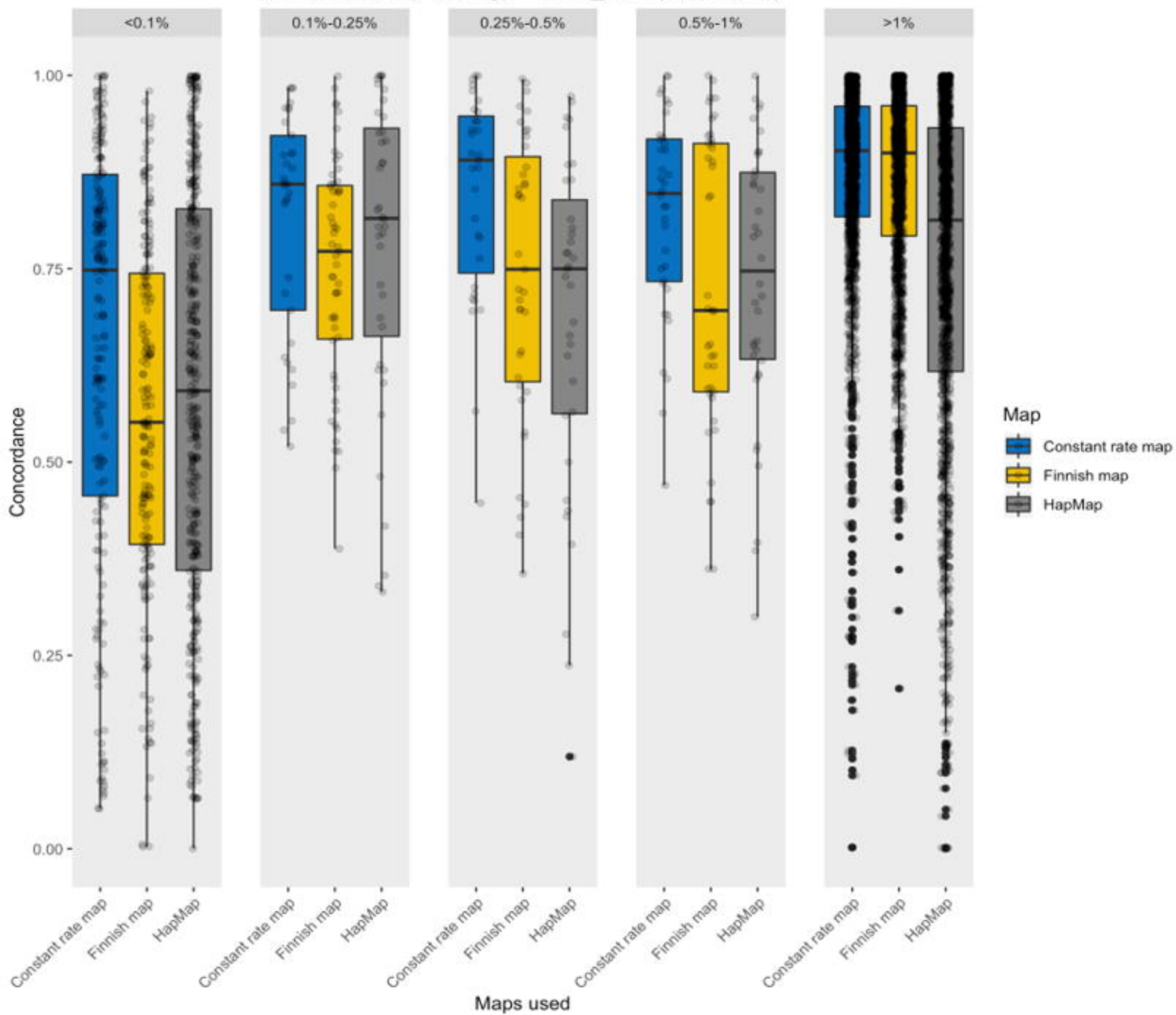




Map  Constant rate map  Finnish map  HapMap



No RefPanels(Phasing) + SISU_2700(Imputation)



With RefPanels(Phasing) + SISU_2700(Imputation)

