

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 Founder population size, demographic changes (eg. population bottlenecks or rapid
3 expansion) can lead to variation in recombination rates across different populations.
4 Previous research has shown that using population-specific reference panels has a
5 significant effect on downstream population genomic analysis like haplotype phasing,
6 genotype imputation and association, especially in the context of population isolates.
7 Here, we developed a high-resolution recombination rate mapping at 10kb and 50kb
8 scale using high-coverage (20-30x) whole-genome sequenced 55 family trios from
9 Finland and compared it to recombination rates of non-Finnish Europeans (NFE). We
10 tested the downstream effects of the population-specific recombination rates in
11 statistical phasing and genotype imputation in Finns as compared to the same analyses
12 performed by using the NFE-based recombination rates. We found that Finnish
13 recombination rates have a moderately high correlation (Spearman's $\rho = 0.67-0.79$) with
14 NFE, although on average (across all autosomal chromosomes), Finnish rates
15 (2.268 ± 0.4209 cM/Mb) are 12-14% lower than NFE (2.641 ± 0.5032 cM/Mb). Finnish
16 recombination map was found to have no significant effect in haplotype phasing
17 accuracy (switch error rates $\sim 2\%$) and average imputation concordance rates (97-98%
18 for common, 92-96% for low frequency and 78-90% for rare variants). Our results
19 suggest that downstream population genomic analyses like haplotype phasing and
20 genotype imputation mostly depend on population-specific contexts like appropriate
21 reference panels and their sample size, but not on population-specific recombination
22 maps or effective population sizes. Currently, available HapMap recombination maps
23 seem robust for population-specific phasing and imputation pipelines, even in the
24 context of relatively isolated populations like Finland.

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

26 **1. Introduction:**

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

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

38 Besides the International HapMap project that focused on capturing common variants
39 and haplotypes in diverse populations, international WGS-based collaborations like
40 1000 Genomes Project, provided genetic variation data for 20 worldwide populations
41 [7]. This led to further refinement of the recombination maps coupled with
42 methodological advances of using coalescent methods for recombination rate [8, 9].

43 With the rise of international collaborative projects, it was realised that founder
44 populations can often have very unique LD patterns [10], subsequently also displaying
45 unique increased genetics-driven health risks [11], suggesting that population-specific
46 reference datasets should be used to leverage the LD patterns to better detect disease
47 variants in downstream genetic analysis [12]. Genomic analysis methods like

48 haplotype phasing and imputing genotypes require recombination maps and other
49 population genetic parameters as input to obtain optimal results [13, 14, 15, 16]
50 In theis study, we set to test this by 1) estimating recombination rates along the genome
51 in Finnish population using ~55 families of whole-genome sequenced (20-30x) Finns,
52 2) comparing these rates to some other European populations, and 3) comparing the
53 effect of using Finnish recombination rate estimates and cosmopolitan estimates in
54 phasing and imputation errors in Finnish samples.

55 **2. Materials & Methods:**

56 **2.1 Datasets used:**

57 ***Finnish Migraine Families Collection***

58 Whole-genome sequenced trios (n = 55) consisting of the parent-offspring combination
59 were drawn from a large Finnish migraine families collection consisting of 1,589
60 families totalling 8,319 individuals [17]. The trios were used for the recombination map
61 construction using LDHAT version 2. The families were collected over 25 years from
62 various headache clinics in Finland (Helsinki, Turku, Jyväskylä, Tampere, Kemi, and
63 Kuopio) and via advertisements in the national migraine patient organisation web page
64 (<https://migreeni.org/>). The families consist of different pedigree sizes from small to
65 large (1-5+ individuals). Of the 8319 individuals, 5317 have a confirmed migraine
66 diagnosis based on the third edition of the established International Classification for
67 Headache Disorders (ICHD-3) criteria [18].

68 ***EUFAM cohort***

69 To check the phasing accuracy of our Finnish recombination map, we used an
70 independently sourced 49 trios from the European Multicenter Study on Familial
71 Dyslipidemias in Patients with Premature Coronary Heart Disease (EUFAM). Finnish

72 familial combined hyperlipidemia (FCH) families were identified from patients initially
73 admitted to hospitals with premature cardiovascular heart disease (CHD) diagnosis who
74 also had elevated levels of total cholesterol (TC), triglycerides (TG) or both in the \geq
75 90th Finnish population percentile. Those families who had at least one additional first-
76 degree relative also affected with hyperlipidemia were also included in the study apart
77 from individuals with elevated levels of TG. [19, 20, 21].

78 ***FINRISK cohort***

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

87 ***FINNISH reference panel cohort***

88 The whole-genome sequenced samples used were obtained from PCR-free methods and
89 PCR-amplified methods, which was followed by sequencing on a Illumina HiSeq X
90 platform with a mean depth of $\sim 30\times$. The obtained reads were then aligned to the
91 GRCh37 (hg19) human reference genome assembly using BWA-MEM. Best practice
92 guidelines from Genome Analysis Toolkit (GATK) were used to process the BAM files
93 and variant calling. Several criteria were used in this stage for sample exclusion:
94 relatedness (identity-by-descent (IBD) > 0.1), sex mismatches, among several others.
95 Furthermore, samples were filtered based on other criteria such as: non-reference

96 variants, singletons, heterozygous/homozygous variants ratio, insertion/deletion ratio
97 for novel indels, insertion/deletion ratio for indels observed in dbSNP, and
98 transition/transversion ratio.

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

107 The final reference dataset used in this study for imputation consisted of high coverage
108 (20-30x) whole-genome sequence-based reference panel of 2690 individuals from the
109 SISu project (Sequencing Initiative Suomi, <http://www.sisuproject.fi/>, [24]).

110 **2.2 Recombination map construction:**

111 Coalescent-based fine-scale recombination map construction [8] is greatly eased by
112 using trios which provide more accurate haplotype phasing resolution [25]. Hence, we
113 used trio data (n=55, 110 independent parents) from the Finnish Migraine Families
114 Cohort described above. These were filtered primarily using VCFtools [26] and custom
115 R scripts. Firstly, sites were thinned with within 15bp of each other such that only one
116 site remained followed by a filtering step of removing variants with a minor allele
117 frequency of <5% [27]. The resultant data were then phased using family-aware
118 method of SHAPEIT [28] using the standard HapMap recombination map [8, 9],
119 which was then split into segments of ~10000 SNPs with a 1000 SNP overhang on each

120 side of the segments. LDhat version 2 was run for 10^7 iterations with a block penalty of
121 5, every 5000 iterations of them of which the first 10% observations were discarded [8,
122 29]. The CEU based maps, used here for comparison, were obtained similarly using
123 LDhat [29].

124 However, LDHat is computationally intensive, and calculations suggest that the 1000
125 Genomes OMNI data set [30] would be too much computationally intensive to
126 complete [31], hence limiting the maximum number of haplotypes which could be
127 used.

128 To overcome this and make the recombination map independent of the underlying
129 methodology, we used a machine learning method implemented in FastEPRR [31, 32].

130 It supports the use of larger sample sizes, than LDHat and the recombination estimates
131 for sample sizes > 50 , yields smaller variance than LDHat based estimates [31]. The
132 method was then applied to each autosome with overlapping sliding windows (*i.e.*,
133 window size, 50 kb and step length, 25 kb) under default settings for diploid organisms.

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

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

142 **2.3 Phasing and imputation accuracy**

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

147 ***2.3.1 Phasing Accuracy***

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

152 For testing the influence of recombination maps on phasing accuracy, we used three
153 different recombination maps: HapMap, fine-scale Finnish recombination map and a
154 constant background recombination rate (1cM/Mb), to phase the 55 offspring
155 haplotypes without using any reference dataset. To check whether reference panels used
156 during haplotype phasing made any impact on the switch error rates, we used the
157 Finnish SISU based reference (n=2690), to check whether the size of the reference
158 panel made any impact on the results in phasing the offspring's haplotypes (Figure 1).

159 The SER in the offspring's phased haplotypes were then calculated by determining the
160 true offspring haplotypes using data from the parents (98 individuals) with a custom
161 script [34].

162 ***2.3.2 Imputation Accuracy***

163 Imputation concordance was used as the metric for calculating the imputation accuracy.
164 For this, we randomly masked FINRISK CoreExome chip data consisting of 10,480
165 individuals [22] from chromosome 20. To test the role of reference panel size in
166 influencing the imputation accuracy in conjunction with varying the population genetics

167 parameters, we imputed the masked dataset with BEAGLE (Browning *et al.*, 2016)
168 using the Finnish reference panel (n = 2690). The concordance was then calculated
169 between the imputed genotypes and the original masked variants. Masking was done by
170 randomly removing ~10% of variants from the chip dataset.

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

175 **3. Results:**

176 **3.1 Finnish recombination map and its comparison to the HapMap recombination
177 map:**

178 The primary aim of our study was to derive a high-resolution genetic recombination
179 map for Finland and use it for comparative tests in commonly used analyses like
180 haplotype phasing and imputation. To derive a population-specific Finnish
181 recombination map, we used the high-coverage WGS data and an average of different
182 estimation methods (LDHat and FastEPRR). We used the Ne value of 10,000 derived
183 from our extensive testing of different Ne values (See supplementary text) to get the
184 per-generation recombination rates. The average recombination rates of Finnish
185 population isolate depicted 12-14% lower values (autosome-wide average 2.268 ± 0.4209
186 cM/Mb) for all chromosomes compared to CEU based maps (2.641 ± 0.5032 cM/Mb)
187 (Figure 2).

188 These differences in average recombination rates are reflected in the correlation values
189 across all chromosomes (Spearman's $\rho \sim 0.67 - 0.79$) between the developed Finnish
190 map and HapMap based one (Figure 2). We also present a direct comparison between

191 the two maps, of the recombination rates at 5Mb scales, which presents a similar visual
192 pattern of rates across the genome (Supplementary Figure 1).

193 **3.2 Effects of the population-specific recombinations map on haplotype phasing**

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

196 We tested the Finnish map, HapMap map and a constant recombination rate map
197 (1cM/Mb) to understand the effects of population-specific maps on downstream
198 genomic analyses. The phasing accuracy was tested under two different conditions:
199 using no additional reference panel and using an population-specific . SISu v2 reference
200 panel (n= 2690) in phasing. We observed that, on average, SER ranged between 1.8-
201 3.7% (Supplementary Figure 2) across the different chromosomes and recombination
202 maps. We found statistically significant differences within both no-reference panel and
203 the Finnish reference panel results (Kruskal Wallis, p-value = 5.3e-10 and 4.7e-10,
204 respectively; Figure 3). The constant recombination map (1cM/Mb) had significantly
205 higher SER values when compared to the Finnish map or the HapMap map (Figure 3)
206 both when no reference panels were used (p-value = 2.9e-11 and 2.6e-09, respectively)
207 and when the Finnish reference panel was used (p-value = 2.9e-11 and 9.5e-13,
208 respectively). The choice of recombination maps mattered more when no reference
209 panel was used (p-value = 0.0046), however when using the Finnish reference panel, the
210 difference in SER was statistically insignificant (p-value = 0.25).

211 **3.3 Effects of the population-specific recombinations map on genotype imputation**

212 Imputation accuracy was similarly tested using the reference panel under three different
213 recombination map settings. We observed that when the imputation target dataset was
214 phased and imputed using the Finnish reference panel (n=2690) irrespective of the

215 population-specific recombination maps, it had a high imputation accuracy (overall
216 concordance rate ~98%, Figure 4) across MAF bins (>0.1%). Though some differences
217 in concordance rates are seen in for rare variants (MAF <0.1%). The concordance rate
218 was lower when the test dataset was phased without reference panels (concordance rate
219 72~77%, Figure 5).

220 **4. Discussion:**

221 Population isolates like Finland, have had a divergent demographic history as compared
222 to the outbred European populations, with a less historic migration, more fluctuating
223 population sizes and higher incidences of bottleneck events and founder effects [35, 36]
224 This unique demographic history then affects different population genetic parameters,
225 like recombination rates [37]. It has been shown previously that using population-
226 specific genomic reference panels augmented the accuracy of imputation accuracy
227 leading to better mapping of diseases specific variants in GWAS [12]. Since
228 recombination rates (in the form of recombination maps), features in much of the
229 downstream genomic analyses' methods like imputation and haplotype phasing [15,
230 34], we wanted to study their effect on downstream analyses.

231 Firstly, we characterised the Finnish recombination map using high-coverage (~30x)
232 whole-genome sequencing (WGS) samples from large SISu v2 reference panel
233 (n=2690). Previously used recombination maps hail from the HapMap and
234 1000Genomes projects which used sparse genotypic datasets or low-depth sequencing
235 samples. This is a first attempt in creating a recombination map for Finland using
236 population-specific WGS samples. We used two different methods in estimating the
237 recombination rates, to achieve accurate estimates with lower variance [29,31]. In
238 addition, we estimated effective population sizes using identity-by-descent (IBD) based

239 methods [15] for both Finnish and CEU based datasets. The obtained recombination
240 map was then used to test their role and importance in two selected downstream
241 genomic analyses – haplotype phasing and imputation concordance. Since the
242 recombination rate determination requires effective population size estimates, we also
243 tested the role of varying effective population size on these two analyses (See
244 Supplementary Text). The extensive testing of N_e yielded the estimate of 10,000
245 originally derived theoretically [38] and most used commonly for humans fits quite
246 rightly for the recombination map.

247 The Finnish recombinational landscape when compared to the HapMap based map,
248 showed, on average, a high degree of correlation across scales (10, 50kb and 5Mb),
249 however, on average, Finnish recombination rates across chromosomes were found to
250 be lower. Such moderate to high correlations (Figure 2) and similar recombinational
251 landscape (Supplementary Figure 1) could be due to high sharing of recombinations in
252 individuals from closely-related populations. The degree of dissimilarity in the
253 population-level differences between Finnish and mainland Europeans in terms of
254 recombination rates could be due to population-specific demographic processes like
255 founder effects, bottleneck events and migration [39], or chromatin structure PRDM9
256 binding locations, for example [40]. And the broad similarity in terms of correlational
257 structure seen here, reflects a shared ancestral origin of Finns and other mainland
258 Europeans [41]. Other studies on population isolates like Iceland [9] have previously
259 found a high degree of correlation with CEU based maps, albeit with substantial
260 differences as seen here. Previous studies [42] have additionally explored the
261 relationship between recombination rate differences between populations and allele
262 frequency differences, with evidence suggesting that the differences between rates show

263 the selection impact in the past 100,000 years since the out-of-Africa movement of
264 humans.

265 As seen in previous studies, much of the downstream genomic analyses like getting
266 more refined GWAS hits or, accurate copy number variants (CNV) imputation, can be
267 highly improved with the addition/use of population-specific datasets [12]. To test this
268 in the context of population-specific recombination maps, we used them to test the
269 haplotype phasing and imputation accuracy and observed that despite large differences
270 in the effective population sizes between populations, it did not affect the tested metrics.

271 One possible explanation for the insignificant effect seen here is that the role of
272 parameters like effective population size and recombination maps is to scale over the
273 haplotypes for efficient coverage of the whole genome. However, when sufficiently
274 large, population-specific genomic reference panels are available with tens of thousands
275 of haplotypic combinations, such scaling over for specific populations, does not yield
276 in substantial improvements. As we showed here, reference panel size could play an
277 important role in the downstream genomic analyses and in most cases, the current
278 practice of using the standard HapMap recombination map can be reasonably used.

279 Another point of interest here is that the use of different N_e parameters during
280 phasing/imputation might be redundant as we observed no change in the accuracy of our
281 estimates on varying the N_e parameters. Similarly, when using population-specific
282 recombination maps, we did not find any tangible benefits in using them over the
283 current standard maps based on the HapMap data.

284 Our study suggests a couple of important points for future studies: (a) varying effective
285 population size for downstream genomic analyses, such as phasing and imputation,
286 might have a relatively small impact, and it might be better to use the default option of

287 the particular software; (b) when available, it is beneficial to use a population-specific
288 genomic reference panel as they increase the accuracy; (c) HapMap can be used for
289 current downstream genomic analyses like haplotype phasing or genotype imputation in
290 European-based populations. And, if need be, can be substituted for using population-
291 specific maps, as the accuracy rates are quite similar to the population-based maps.
292 Though the sample used here is from a disease cohort but is nevertheless representative
293 of Finland's population and hence provides a reasonable recombination rate estimates.
294 On the other hand, our reliance on disease cohorts could lead to minor variation in the
295 resultant recombination. Though as we share a similar out-of-Africa origin, much of our
296 history is shared and though biological differences in the recombinational landscape do
297 exist between different populations, much of the downstream genomic analyses
298 (haplotyping, imputation or, GWAS), might not be affected by recombination map or
299 values of effective population size.

300 **Funding**

301 This work was financially supported by the Academy of Finland (251217 and 255847 to
302 S.R.). S.R. was further supported by the Academy of Finland, Center of Excellence for
303 Complex Disease Genetics, the Finnish Foundation for Cardiovascular Research,
304 Biocentrum Helsinki, and the Sigrid Jusélius Foundation. S.H. was supported by
305 FIMM-EMBL PhD program doctoral funding and I.S. by Academy of Finland
306 Postdoctoral Fellowship (298149). V.S. was supported by the Finnish Foundation for
307 Cardiovascular Research. T.T. was supported by Academy of Finland grant number
308 315589.

309 **Acknowledgements**

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

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

317

318

319 **References**

- 320 1. Baudat F, Buard J, Grey C, Fledel-Alon A, Ober C, Przeworski M et al. PRDM9
321 is a major determinant of meiotic recombination hotspots in humans and mice.
322 Science 2009; 327: 836–840
- 323 2. Daly MJ, Rioux JD, Schaffner SF, Hudson TJ, Lander ES. High-resolution
324 haplotype structure in the human genome. Nature Genetics 2001; 29: 229–232
- 325 3. Hudson RR, Kaplan NL. Statistical properties of the number of recombination
326 events in the history of a sample of DNA sequences. Genetics 1985; 111: 147–
327 164
- 328 4. Chan AH, Jenkins PA, Song YS. Genome-wide fine-scale recombination rate
329 variation in *Drosophila melanogaster*. PLoS Genet 2012; 8: e1003090
- 330 5. McVean GA, Myers SR, Hunt S, Deloukas P, Bentley DR, Donnelly P. The
331 fine-scale structure of recombination rate variation in the human genome.
332 Science 2004; 304: 581-584

333 6. Myers S, Bottolo L, Freeman C, McVean G, Donnelly P. A fine-scale map of
334 recombination rates and hotspots across the human genome. *Science* 2005; 310:
335 321-324.

336 7. Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO et al. A
337 global reference for human genetic variation. *Nature* 2015; 526: 68-74

338 8. Auton A, McVean G. Recombination rate estimation in the presence of hotspots.
339 *Genome Res* 2007; 17: 1219-1227.

340 9. Kong A, Thorleifsson G, Gudbjartsson DF, Masson G, Sigurdsson A,
341 Jonasdottir A et al. Fine-scale recombination rate differences between sexes,
342 populations and individuals. *Nature* 2010; 467: 1099-1103.

343 10. Service S, DeYoung J, Karayiorgou M, Roos JL, Pretorius H, Bedoya G et al.
344 Magnitude and distribution of linkage disequilibrium in population isolates and
345 implications for genome-wide association studies. *Nat Genet* 2006; 38: 556-560.

346 11. Peltonen L, Jalanko A, Varilo T. Molecular genetics of the Finnish disease
347 heritage. *Hum Mol Genet* 1999; 8: 1913-1923.

348 12. Surakka I, Kristiansson K, Anttila V, Inouye M, Barnes C, Moutsianas L et al.
349 Founder population-specific HapMap panel increases power in GWA studies
350 through improved imputation accuracy and CNV tagging. *Genome Res* 2010;
351 20: 1344-1351.

352 13. Tewhey R, Bansal V, Torkamani A, Topol EJ, Schork NJ. The importance of
353 phase information for human genomics. *Nat Rev Genet* 2011; 12: 215-223.

354 14. Browning SR, Browning BL. Haplotype phasing: existing methods and new
355 developments. *Nat Rev Genet* 2011; 12: 703-714.

356 15. Browning BL, Browning SR. Genotype Imputation with Millions of Reference
357 Samples. *Am J Hum Genet* 2016; 98: 116-126.

358 16. Delaneau O, Zagury JF, Marchini J. Improved whole-chromosome phasing for
359 disease and population genetic studies. *Nat Methods* 2013; 10: 5-6.

360 17. Gormley P, Kurki MI, Hiekkala ME, Veerapen K, Häppölä P, Mitchell AA et al.
361 Common Variant Burden Contributes to the Familial Aggregation of Migraine
362 in 1,589 Families. *Neuron* 2018; 98: 743-753.e4.

363 18. The International Classification of Headache Disorders, 3rd edition (beta
364 version). *Cephalgia* 2013; 33: 629-808.

365 19. Borodulin K, Vartiainen E, Peltonen M, Jousilahti P, Juolevi A, Laatikainen T et
366 al. Forty-year trends in cardiovascular risk factors in Finland. *Eur J Public
367 Health* 2015; 25: 539-546.

368 20. Porkka KV, Nuotio I, Pajukanta P, Ehnholm C, Suurinkeroinen L, Syvänne M et
369 al. Phenotype expression in familial combined hyperlipidemia. *Atherosclerosis*
370 1997; 133: 245-253.

371 21. Ripatti P, Rämö JT, Söderlund S, Surakka I, Matikainen N, Pirinen M et al. The
372 Contribution of GWAS Loci in Familial Dyslipidemias. *PLOS Genetics* 2016;
373 12: e1006078.

374 22. Vartiainen E, Laatikainen T, Peltonen M, Juolevi A, Mannisto S, Sundvall J et
375 al. Thirty-five-year trends in cardiovascular risk factors in Finland. *International
376 Journal of Epidemiology* 2009; 39: 504–518.

377 23. Li H. Toward better understanding of artifacts in variant calling from high-
378 coverage samples. *Bioinformatics* 2014; 30: 2843–2851.

379 24. Mart Kals, Tiit Nikopensius, Kristi Läll, Kalle Pärn, Timo Tõnis Sikka, Jaana
380 Suvisaari, Veikko Salomaa, Samuli Ripatti, Aarno Palotie, Andres Metspalu,
381 Tõnu Esko, Priit Palta, Reedik Mägi Advantages of genotype imputation with
382 ethnically matched reference panel for rare variant association analyses bioRxiv
383 579201; doi: <https://doi.org/10.1101/579201>

384 25. Roach JC, Glusman G, Hubley R, Montsaroff SZ, Holloway AK, Mauldin DE et
385 al. Chromosomal haplotypes by genetic phasing of human families. Am J Hum
386 Genet 2011; 89: 382-397.

387 26. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA et al. The
388 variant call format and VCFtools. Bioinformatics 2011; 27: 2156-2158.

389 27. Steviston LS, Woerner AE, Kidd JM, Kelley JL, Veeramah KR, McManus KF et
390 al. The Time Scale of Recombination Rate Evolution in Great Apes. Mol Biol
391 Evol 2016; 33: 928-945.

392 28. O'Connell J, Gurdasani D, Delaneau O, Pirastu N, Ulivi S, Cocca M et al. A
393 general approach for haplotype phasing across the full spectrum of relatedness.
394 PLoS Genet 2014; 10: e1004234.

395 29. Auton A, Fledel-Alon A, Pfeifer S, Venn O, Séguirel L, Street T et al. A fine-
396 scale chimpanzee genetic map from population sequencing. Science 2012; 336:
397 193-198.

398 30. Abecasis GR, Auton A, Brooks LD, DePristo MA, Durbin RM, Handsaker RE
399 et al. An integrated map of genetic variation from 1,092 human genomes. Nature
400 2012; 491: 56-65.

401 31. Gao F, Ming C, Hu W, Li H. New Software for the Fast Estimation of
402 Population Recombination Rates (FastEPRR) in the Genomic Era. *G3*
403 (Bethesda) 2016; 6: 1563-1571.

404 32. Lin K, Futschik A, Li H. A fast estimate for the population recombination rate
405 based on regression. *Genetics* 2013; 194: 473-484.

406 33. Bansal V. Integrating read-based and population-based phasing for dense and
407 accurate haplotyping of individual genomes. *Bioinformatics* 2019; 35: i242-
408 i248.

409 34. Loh PR, Danecek P, Palamara PF, Fuchsberger C, A Reshef Y, K Finucane H et
410 al. Reference-based phasing using the Haplotype Reference Consortium panel.
411 *Nat Genet* 2016; 48: 1443-1448.

412 35. Martin AR, Karczewski KJ, Kerminen S, Kurki MI, Sarin AP, Artomov M et al.
413 Haplotype Sharing Provides Insights into Fine-Scale Population History and
414 Disease in Finland. *Am J Hum Genet* 2018; 102: 760-775.

415 36. Kerminen S, Havulinna AS, Hellenthal G, Martin AR, Sarin AP, Perola M et al.
416 Fine-Scale Genetic Structure in Finland. *G3* (Bethesda) 2017; 7: 3459-3468.

417 37. Wang J, Santiago E, Caballero A. Prediction and estimation of effective
418 population size. *Heredity (Edinb)* 2016; 117: 193-206.

419 38. Takahata N, Satta Y, Klein J. Divergence time and population size in the lineage
420 leading to modern humans. *Theor Popul Biol* 1995; 48: 198-221.

421 39. Novembre J, Johnson T, Bryc K, Kutalik Z, Boyko AR, Auton A et al. Genes
422 mirror geography within Europe. *Nature* 2008; 456: 98-101.

423 40. Ségurel L. The complex binding of PRDM9. *Genome Biol* 2013; 14: 112.

424 41. Mallick S, Li H, Lipson M, Mathieson I, Gymrek M, Racimo F et al. The
425 Simons Genome Diversity Project: 300 genomes from 142 diverse populations.
426 Nature 2016; 538: 201-206.

427 42. Keinan A, Reich D. Human population differentiation is strongly correlated with
428 local recombination rate. PLoS Genet 2010; 6: e1000886.

429

430

431

432

433

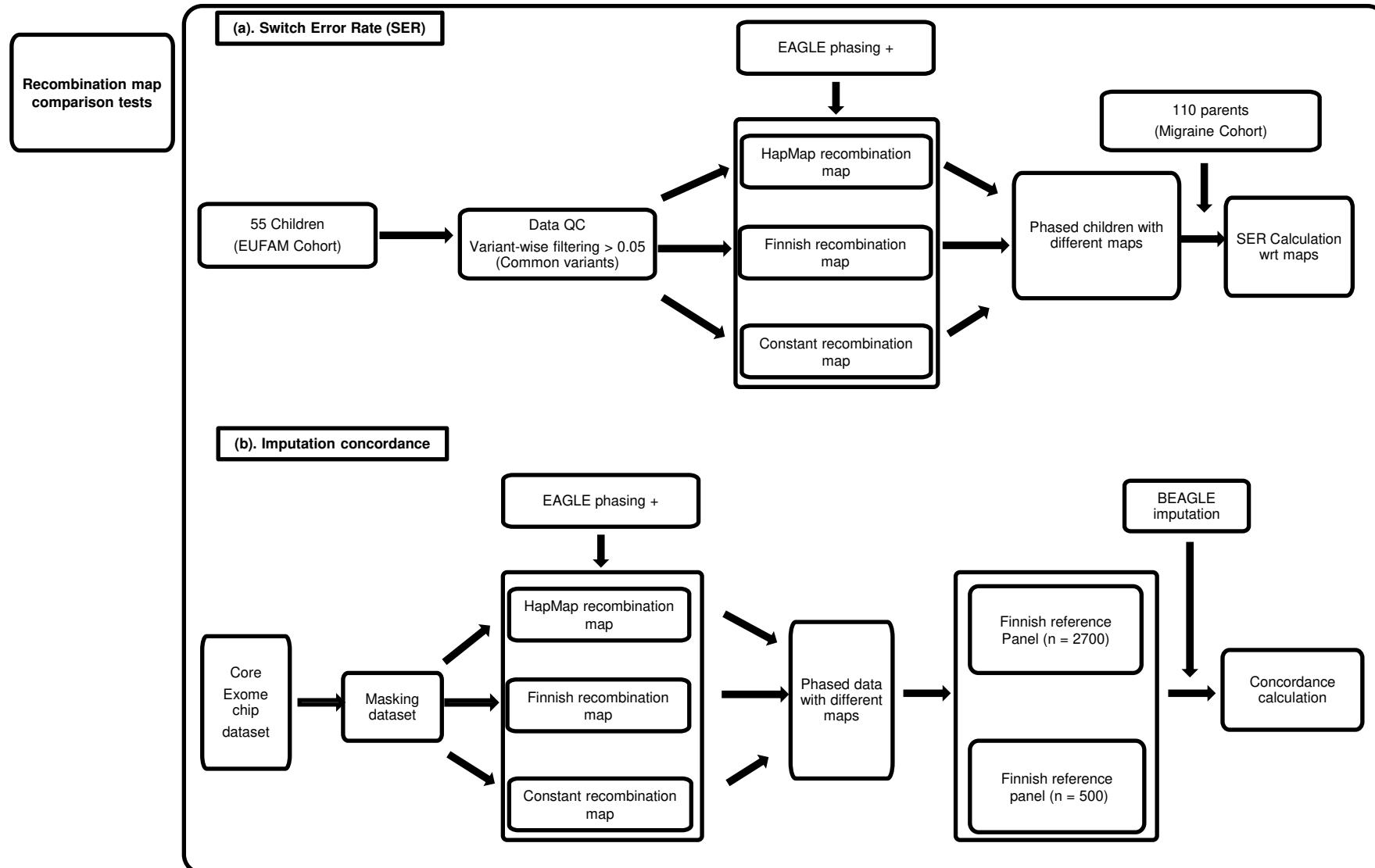
434

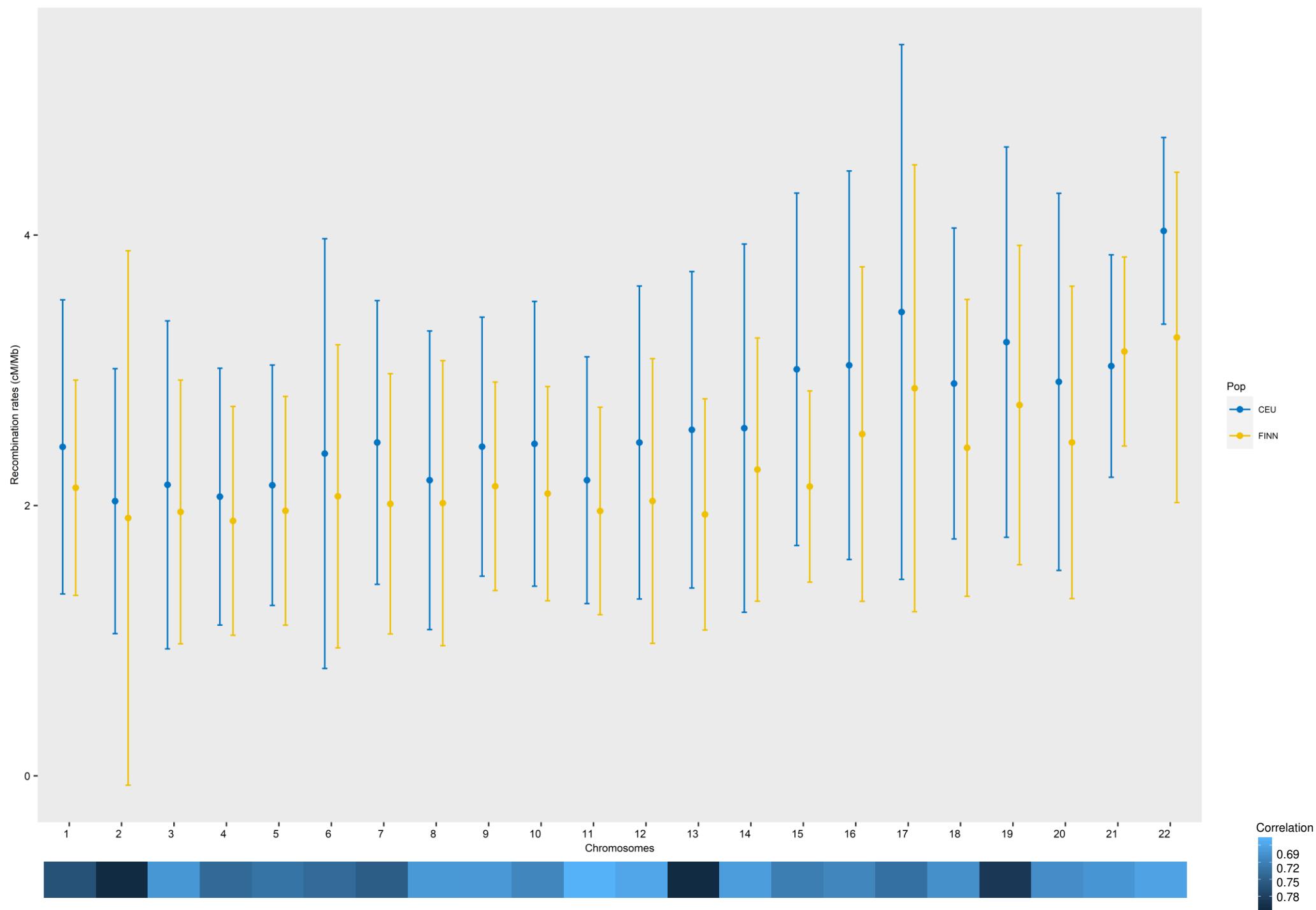
435

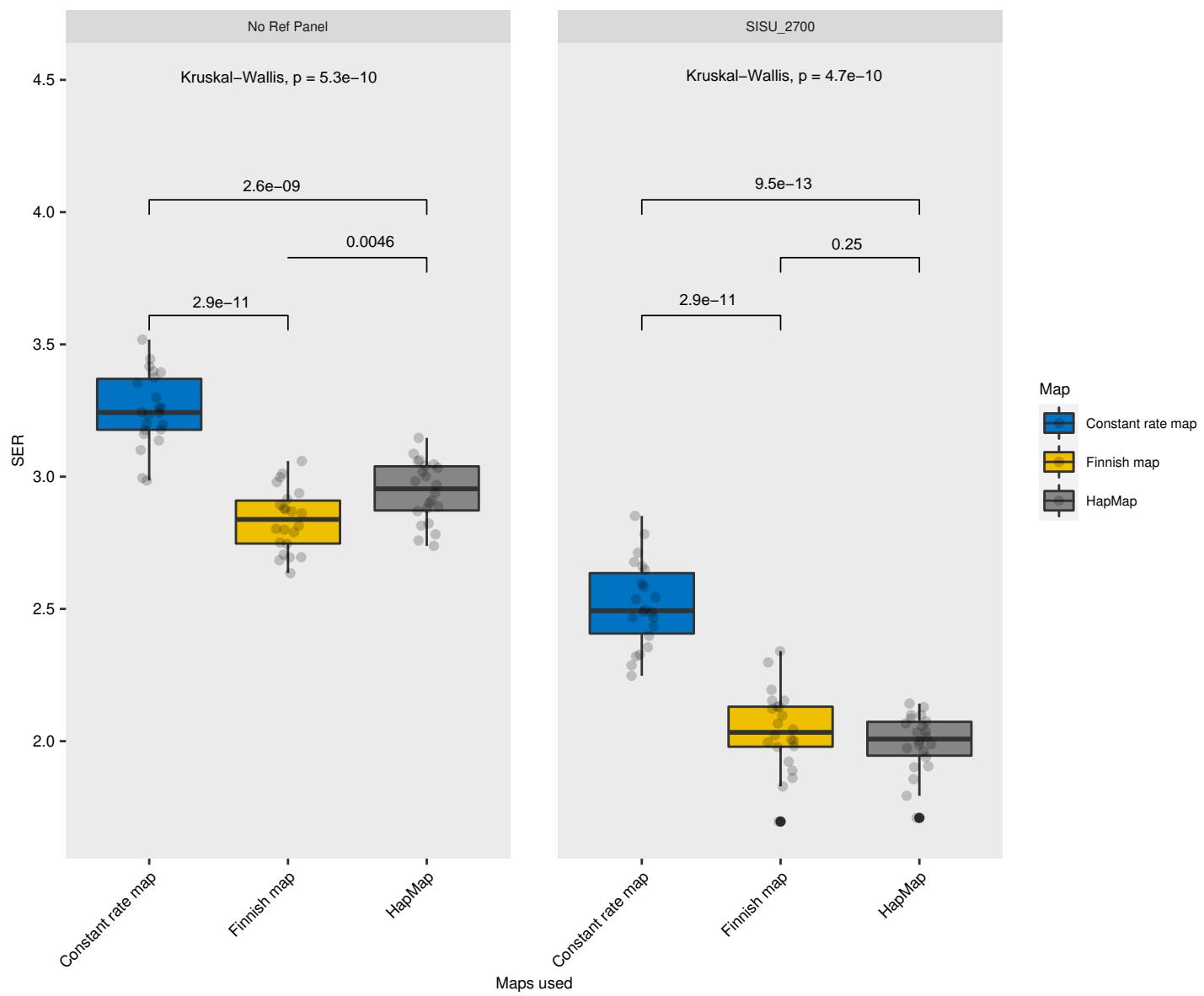
436 **Figure 1:** Flowchart overview of the analyses and comparisons performed.

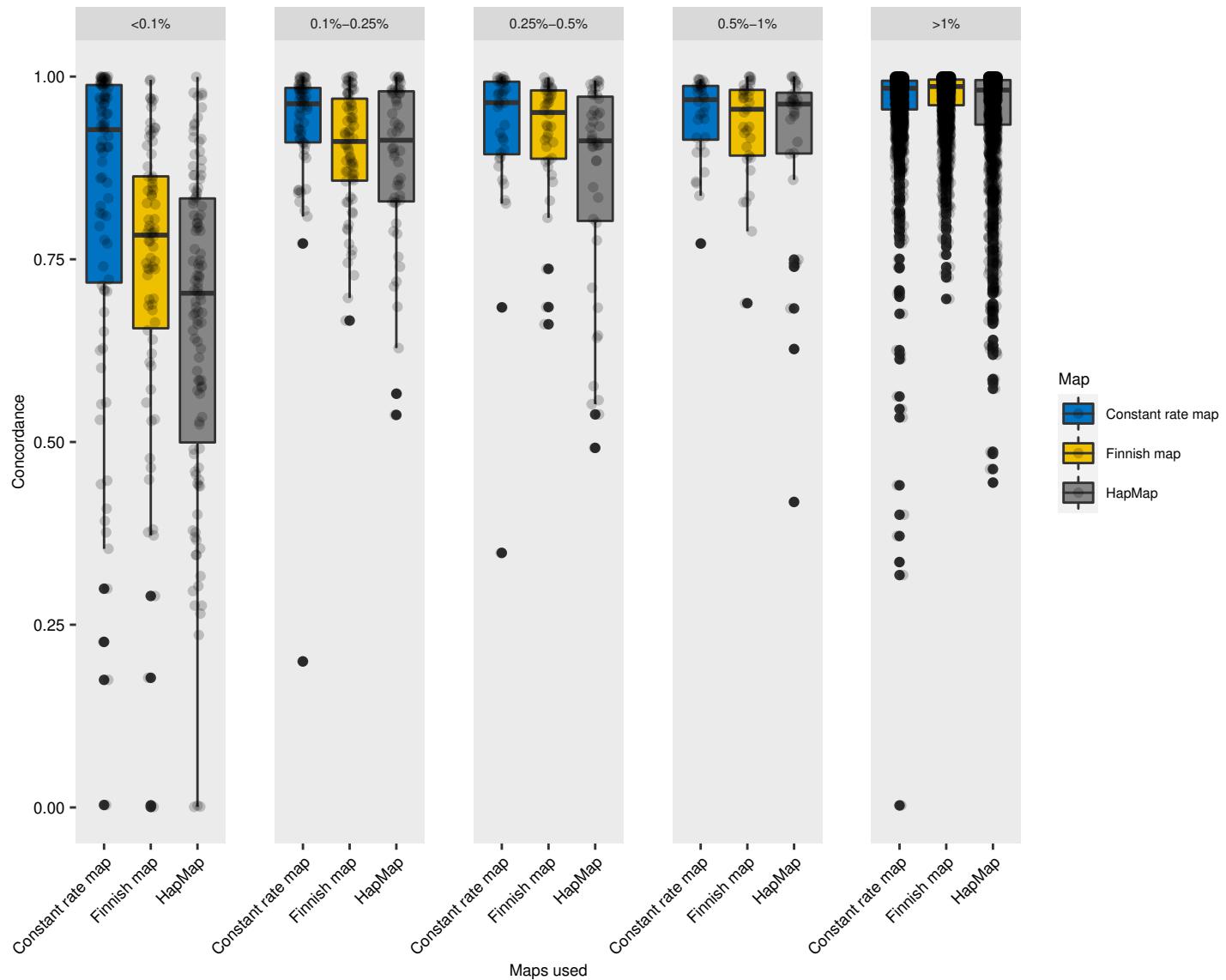
437 **Figure 2:** Average (\pm standard deviation) recombination rates of Finnish v/s CEU per
438 autosome measured in cM/Mb and Correlation between Finnish and CEU
439 recombination rates across all chromosomes. The comparisons are made for similar
440 physical positions.

441 **Figure 3:** Statistical comparison of Switch Error Rates across all autosomes calculated
442 for all children in the trios using different recombination maps with respect to different
443 reference panel conditions (absent or present). The p-values are shown at the top of each
444 panel from Kruskal Wallis ANOVA testing between panel groups and ones between
445 boxplots for within-group comparisons.









No RefPanels(Phasing) + SISU_2700(Imputation)

